

日 本 国 特 許 庁
JAPAN PATENT OFFICE

31.03.03

別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office

出 願 年 月 日
Date of Application:

2002年12月 4日

出 願 番 号
Application Number:

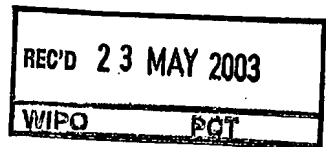
特願2002-352645

[ST.10/C]:

[JP2002-352645]

出 願 人
Applicant(s):

石原産業株式会社

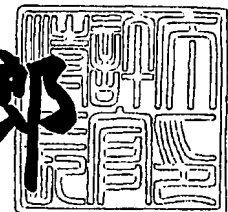


**PRIORITY
DOCUMENT**
SUBMITTED OR TRANSMITTED IN
COMPLIANCE WITH RULE 17.1(a) OR (b)

2003年 5月 9日

特 許 庁 長 官
Commissioner,
Japan Patent Office

太田信一郎



出証番号 出証特2003-3033963

【書類名】 特許願

【整理番号】 186699

【提出日】 平成14年12月 4日

【あて先】 特許庁長官殿

【国際特許分類】 G06F 17/00

【発明者】

 【住所又は居所】 滋賀県草津市西渋川二丁目3番1号 石原産業株式会社
 中央研究所内

 【氏名】 石川 俊夫

【発明者】

 【住所又は居所】 滋賀県草津市西渋川二丁目3番1号 石原産業株式会社
 中央研究所内

 【氏名】 久米 隆志

【特許出願人】

 【識別番号】 000000354

 【住所又は居所】 大阪府大阪市西区江戸堀一丁目3番15号

 【氏名又は名称】 石原産業株式会社

【代理人】

 【識別番号】 100062144

 【弁理士】

 【氏名又は名称】 青山 葆

【選任した代理人】

 【識別番号】 100086405

 【弁理士】

 【氏名又は名称】 河宮 治

【選任した代理人】

 【識別番号】 100098280

 【弁理士】

 【氏名又は名称】 石野 正弘

【先の出願に基づく優先権主張】

【出願番号】 特願2002-102743

【出願日】 平成14年 4月 4日

【手数料の表示】

【予納台帳番号】 013262

【納付金額】 21,000円

【提出物件の目録】

【物件名】 明細書 1

【物件名】 図面 1

【物件名】 要約書 1

【物件名】 委任状 1

【提出物件の特記事項】 手続補足書により提出する。

【ブルーフの要否】 要

【書類名】 明細書

【発明の名称】 データ解析装置および方法

【特許請求の範囲】

【請求項1】 生体の状態または時間とともに確率的に発生する生体の状態の変化を目的変数とし、複数の遺伝子発現の量および／または細胞内物質の量を説明変数とする相関モデルを決定するデータ解析装置であって、

生体の状態或いはそれを導出するデータまたは時間とともに確率的に発生する生体の状態の変化に関するデータと、複数の遺伝子発現の量および／または細胞内物質の量からなるサンプルの集合を入力する入力手段と、

(1)説明変数を選択する選択手段と、

(2)部分最小自乗法を実行して交差検証成績を計算する計算手段または前記生体の状態の変化に関するデータにカプラン・マイヤー法又はカトラー・エデラー法による生命表を適用して変化の発生しなかったものの確率を計算して得られた確率を、仮定した分布に基づいた変換または仮定を前提としない変換をし、該変換結果を目的変数とする部分最小自乗法を実行して交差検証成績を計算する計算手段と、

(3)前記(2)の計算手段の結果を評価し、説明変数の採用、不採用を判定する評価判定手段とを有し、

(4)前記(1)の選択手段と前記(2)の計算手段と前記(3)の評価判定手段とを実行して部分最小自乗法モデルの少なくとも交差検証成績を独立変数として持つ関数を改善し続けて部分最小自乗法モデルを決定する決定手段とからなることを特徴とするデータ解析装置。

【請求項2】 目的変数が生体の状態であって、前記入力手段で入力するデータが生体の状態或いはそれを導出するデータであって、前記(2)の計算手段が部分最小自乗法を実行して交差検証成績を計算する計算手段であることを特徴とする請求項1に記載のデータ解析装置。

【請求項3】 目的変数が時間とともに確率的に発生する生体の状態の変化であって、前記入力手段で入力するデータが時間とともに確率的に発生する生体の状態の変化に関するデータであって、前記(2)の計算手段が前記生体の状態の

変化に関するデータに Kaplan-Meier 法又は Kotlér-Eder 法による生命表を適用して変化の発生しなかったものの確率を計算して得られた確率を、仮定した分布に基づいた変換または仮定を前提としない変換をし、該変換結果を目的変数とする部分最小自乗法を実行して交差検証成績を計算する計算手段であることを特徴とする請求項 1 に記載のデータ解析装置。

【請求項 4】 さらに、前記の決定手段にて決定された部分最小自乗法モデルに採用されている説明変数を用い、統計的手法又は多変量解析手法によるモデルを構築する最終モデル決定手段を備えることを特徴とする請求項 1、2 又は 3 に記載のデータ解析装置。

【請求項 5】 前記の選択手段において、説明変数を逐次取捨選択することを特徴とする請求項 1～4 のいずれかに記載のデータ解析装置。

【請求項 6】 前記の選択手段において、遺伝的アルゴリズムを用いて説明変数を選択することを特徴とする請求項 1～4 のいずれかに記載のデータ解析装置。

【請求項 7】 前記の計算手段において、1 個のサンプルを逐次除外して部分最小自乗法を実行して交差検証成績を計算することを特徴とする請求項 1～6 のいずれかに記載のデータ解析装置。

【請求項 8】 前記の計算手段において、複数のサンプルを逐次除外して部分最小自乗法を実行して交差検証成績を計算することを特徴とする請求項 1～6 のいずれかに記載のデータ解析装置。

【請求項 9】 前記計算手段において、各計算において除外したサンプルの遺伝子発現から予測される生体の状態を示す目的変数値と、前記除外したサンプルの生体の状態を示す目的変数値との誤差の代表値を求め、交差検証成績の指標として当該誤差を用いることを特徴とする請求項 7 又は 8 に記載されたデータ解析装置。

【請求項 10】 前記関数が交差検証成績であることを特徴とする請求項 1～9 のいずれかに記載のデータ解析装置。

【請求項 11】 前記関数が交差検証成績と選択された説明変数の個数との関数であることを特徴とする請求項 1～9 のいずれかに記載のデータ解析装置。

【請求項 1 2】 前記の決定手段において、少なくとも交差検証成績を独立変数として持つ関数を改善しながら評価判定を繰り返すことを特徴とする請求項 5 に記載されたデータ解析装置。

【請求項 1 3】 前記 (1) の選択手段と前記 (2) の計算手段とを複数のコンピュータで実行させることを特徴とする請求項 1 ～ 1 2 のいずれかに記載のデータ解析装置。

【請求項 1 4】 請求項 1、2、3 又は 4 で決定された相関モデル及び予測対象のサンプルについて当該モデルにおいて採用された説明変数を入力する入力手段と、入力された該説明変数に基づいて該サンプルの生体の状態を予測判定する予測判定手段からなることを特徴とするデータ解析装置。

【請求項 1 5】 生体の状態を名義尺度、順序尺度或いは連続量で表現する目的変数とする請求項 2 に記載データ解析装置。

【請求項 1 6】 最終モデル決定手段が用いる前記の統計的手法又は多変量解析手法が、比例ハザード法又はパラメトリックな分布にあてはめた回帰分析法であることを特徴とする請求項 2 又は 4 に記載のデータ解析装置。

【請求項 1 7】 生体の状態または時間とともに確率的に発生する生体の状態の変化を目的変数とし、複数の遺伝子発現の量および／または細胞内物質の量を説明変数とする相関モデルを決定するデータ解析方法であって、

生体の状態或いはそれを導出するデータまたは時間とともに確率的に発生する生体の状態の変化に関するデータと、複数の遺伝子発現の量および／または細胞内物質の量からなるサンプルの集合を入力する入力ステップと、

(1) 説明変数を選択する選択ステップと、

(2) 部分最小自乗法を実行して交差検証成績を計算する計算ステップまたは前記生体の状態の変化に関するデータにカプラン・マイヤー法又はカトラー・エデラー法による生命表を適用して変化の発生しなかったものの確率を計算して得られた確率を、仮定した分布に基づいた変換または仮定を前提としない変換をし、該変換結果を目的変数とする部分最小自乗法を実行して交差検証成績を計算する計算ステップと、

(3) 前記 (2) の計算ステップの結果を評価し、説明変数の採用、不採用を判定す

る評価判定ステップとを有し、

(4)前記(1)の選択ステップと前記(2)の計算ステップと前記(3)の評価判定ステップとを実行して部分最小自乗法モデルの少なくとも交差検証成績を独立変数として持つ関数を改善し続けて部分最小自乗法モデルを決定する決定ステップとからなることを特徴とするデータ解析方法。

【請求項18】 目的変数が生体の状態であって、前記入力ステップで入力するデータが生体の状態或いはそれを導出するデータであって、前記(2)の計算ステップが部分最小自乗法を実行して交差検証成績を計算する計算ステップであることを特徴とする請求項17に記載のデータ解析方法。

【請求項19】 目的変数が時間とともに確率的に発生する生体の状態の変化であって、前記入力ステップで入力するデータが時間とともに確率的に発生する生体の状態の変化に関するデータであって、前記(2)の計算ステップが前記生体の状態の変化に関するデータにカプラン・マイヤー法又はカトラー・エデラー法による生命表を適用して変化の発生しなかったものの確率を計算して得られた確率を、仮定した分布に基づいた変換または仮定を前提としない変換をし、該変換結果を目的変数とする部分最小自乗法を実行して交差検証成績を計算する計算ステップであることを特徴とする請求項17に記載のデータ解析方法。

【請求項20】 さらに、前記の決定ステップにて決定された部分最小自乗法モデルに採用されている説明変数を用い、統計的手法又は多変量解析手法によるモデルを構築する最終モデル決定ステップを備えることを特徴とする請求項17、18又は19に記載のデータ解析方法。

【請求項21】 前記の選択ステップにおいて、説明変数を逐次取捨選択することを特徴とする請求項17～20のいずれかに記載のデータ解析方法。

【請求項22】 前記の選択ステップにおいて、遺伝的アルゴリズムを用いて説明変数を選択することを特徴とする請求項17～20のいずれかに記載のデータ解析方法。

【請求項23】 前記の計算ステップにおいて、1個のサンプルを逐次除外して部分最小自乗法を実行して交差検証成績を計算することを特徴とする請求項17～22のいずれかに記載のデータ解析方法。

【請求項 24】 前記の計算ステップにおいて、複数のサンプルを逐次除外して部分最小自乗法を実行して交差検証成績を計算することを特徴とする請求項 17～22 のいずれかに記載のデータ解析方法。

【請求項 25】 前記計算ステップにおいて、各計算において除外したサンプルの遺伝子発現から予測される生体の状態を示す目的変数値と、前記除外したサンプルの生体の状態を示す目的変数値との誤差の代表値を求め、交差検証成績の指標として当該誤差を用いることを特徴とする請求項 23 又は 24 に記載されたデータ解析方法。

【請求項 26】 前記関数が交差検証成績であることを特徴とする請求項 17～25 のいずれかに記載のデータ解析方法。

【請求項 27】 前記関数が交差検証成績と選択された説明変数の個数との関数であることを特徴とする請求項 17～25 のいずれかに記載のデータ解析方法。

【請求項 28】 前記決定ステップにおいて、少なくとも交差検証成績を独立変数として持つ関数を改善しながら評価判定を繰り返すことを特徴とする請求項 21 に記載されたデータ解析方法。

【請求項 29】 前記 (1) の選択ステップと前記 (2) の計算ステップとを複数のコンピュータで実行させることを特徴とする請求項 17～28 のいずれかに記載のデータ解析方法。

【請求項 30】 請求項 17、18、19 又は 20 で決定された相関モデル及び予測対象のサンプルについて当該モデルにおいて採用された説明変数を入力する入力ステップと、入力された該説明変数に基づいて該サンプルの生体の状態を予測判定する予測判定ステップからなることを特徴とするデータ解析方法。

【請求項 31】 生体の状態を名義尺度、順序尺度或いは連続量で表現する目的変数とする請求項 18 に記載データ解析方法。

【請求項 32】 前記の統計的手法又は多変量解析手法が、比例ハザード法又はパラメトリックな分布にあてはめた回帰分析法によるモデルを構築する最終モデル決定ステップとからなることを特徴とする請求項 18 又は 20 に記載のデータ解析方法。

【請求項33】 生体の状態または時間とともに確率的に発生する生体の状態の変化を目的変数とし、複数の遺伝子発現の量および／または細胞内物質の量を説明変数とする相関モデルを決定する、コンピュータにより実行されるデータ解析プログラムであって、

生体の状態或いはそれを導出するデータまたは時間とともに確率的に発生する生体の状態の変化に関するデータと、複数の遺伝子発現の量および／または細胞内物質の量からなるサンプルの集合を入力する入力ステップと、

(1)説明変数を選択する選択ステップと、

(2)部分最小自乗法を実行して交差検証成績を計算する計算ステップまたは前記生体の状態の変化に関するデータにカプラン・マイヤー法又はカトラー・エデラー法による生命表を適用して変化の発生しなかったものの確率を計算して得られた確率を、仮定した分布に基づいた変換または仮定を前提としない変換をし、該変換結果を目的変数とする部分最小自乗法を実行して交差検証成績を計算する計算ステップと、

(3)前記(2)の計算ステップの結果を評価し、説明変数の採用、不採用を判定する評価判定ステップとを有し、

(4)前記(1)の選択ステップと前記(2)の計算ステップと前記(3)の評価判定ステップとを実行して部分最小自乗法モデルの少なくとも交差検証成績を独立変数として持つ関数を改善し続けて部分最小自乗法モデルを決定する決定ステップとからなることを特徴とするデータ解析プログラム。

【請求項34】 目的変数が生体の状態であって、前記入力ステップで入力するデータが生体の状態或いはそれを導出するデータであって、前記(2)の計算ステップが部分最小自乗法を実行して交差検証成績を計算する計算ステップであることを特徴とする請求項33に記載のデータ解析プログラム。

【請求項35】 目的変数が時間とともに確率的に発生する生体の状態の変化であって、前記入力ステップで入力するデータが時間とともに確率的に発生する生体の状態の変化に関するデータであって、前記(2)の計算ステップが前記生体の状態の変化に関するデータにカプラン・マイヤー法又はカトラー・エデラー法による生命表を適用して変化の発生しなかったものの確率を計算して得られた

確率を、仮定した分布に基づいた変換または仮定を前提としない変換をし、該変換結果を目的変数とする部分最小自乗法を実行して交差検証成績を計算する計算ステップであることを特徴とする請求項33に記載のデータ解析プログラム。

【請求項36】 さらに、前記の決定ステップにて決定された部分最小自乗法モデルに採用されている説明変数を用い、統計的手法又は多変量解析手法によるモデルを構築する最終モデル決定ステップを備えることを特徴とする請求項33、34又は35に記載のデータ解析プログラム。

【請求項37】 前記の選択ステップにおいて、説明変数を逐次取捨選択することを特徴とする請求項33～36のいずれかに記載のデータ解析プログラム。

【請求項38】 前記の選択ステップにおいて、遺伝的アルゴリズムを用いて説明変数を選択することを特徴とする請求項33～36のいずれかに記載のデータ解析プログラム。

【請求項39】 前記の計算ステップにおいて、1個のサンプルを逐次除外して部分最小自乗法を実行して交差検証成績を計算することを特徴とする請求項33～38のいずれかに記載のデータ解析プログラム。

【請求項40】 前記の計算ステップにおいて、複数のサンプルを逐次除外して部分最小自乗法を実行して交差検証成績を計算することを特徴とする請求項33～38のいずれかに記載のデータ解析プログラム。

【請求項41】 前記計算ステップにおいて、各計算において除外したサンプルの遺伝子発現から予測される生体の状態を示す目的変数値と、前記除外したサンプルの生体の状態を示す目的変数値との誤差の代表値を求め、交差検証成績の指標として当該誤差を用いることを特徴とする請求項39又は40に記載されたデータ解析プログラム。

【請求項42】 前記関数が交差検証成績であることを特徴とする請求項33～41のいずれかに記載のデータ解析プログラム。

【請求項43】 前記関数が交差検証成績と選択された説明変数の個数との関数であることを特徴とする請求項33～41のいずれかに記載のデータ解析プログラム。

【請求項44】 前記決定ステップにおいて、少なくとも交差検証成績を独立変数として持つ関数を改善しながら評価判定を繰り返すことを特徴とする請求項37に記載されたデータ解析プログラム。

【請求項45】 前記(1)の選択ステップと前記(2)の計算ステップとを複数のコンピュータで実行させることを特徴とする請求項33～44のいずれかに記載のデータ解析プログラム。

【請求項46】 請求項33、34、35又は36で決定された相関モデル及び予測対象のサンプルについて当該モデルにおいて採用された説明変数を入力する入力ステップと、入力された該説明変数に基づいて該サンプルの生体の状態を予測判定する予測判定ステップからなることを特徴とするデータ解析プログラム。

【請求項47】 生体の状態を名義尺度、順序尺度或いは連続量で表現する目的変数とする請求項34に記載データ解析プログラム。

【請求項48】 前記の統計的手法又は多変量解析手法が、比例ハザード法又はパラメトリックな分布にあてはめた回帰分析法によるモデルを構築する最終モデル決定ステップとからなることを特徴とする請求項34又は36に記載のデータ解析プログラム。

【請求項49】 上記の説明変数の選択において、初期状態では説明変数を全く含まないことを特徴とする請求項37に記載されたプログラム。

【請求項50】 上記の説明変数の選択において、初期状態では全説明変数を含むことを特徴とする請求項37に記載されたプログラム。

【請求項51】 上記の生体の状態が病気のタイプをあらわす測定値、病気の重篤度をあらわす測定値、病気のタイプをあらわす医療診断の結果、病気の重篤度をあらわす医療診断の結果、あるいはそれらを2次加工した数値であることを特徴とする請求項37から50のいずれかに記載されたプログラム。

【請求項52】 請求項33から請求項48のいずれかに記載されたプログラムを記録した、コンピュータにより読み取り可能な記録媒体。

【請求項53】 実質的にジーンバンクアクセッション番号がU15085、M23452、X52479、U70426、H57330及びS69790からなる遺伝子群の発現を検出すること

を特徴とするびまん性大細胞型Bリンパ腫の重篤度検定用の細胞内物質測定機材および測定方法並びにびまん性大細胞型Bリンパ腫の重篤度検定方法。

【請求項54】 さらにジーンバンクアクセッション番号がU03398、M65066、AK001546、BC003536、X00437、U12979、H96306、AA830781及びAA804793からなる群から選択される少なくとも一つの遺伝子の発現を検出することを特徴とする請求項53に記載のびまん性大細胞型Bリンパ腫の重篤度検定用の細胞内物質測定機材および測定方法並びにびまん性大細胞型Bリンパ腫の重篤度検定方法。

【請求項55】 実質的にジーンバンクアクセッション番号がAA598572、AA703058及びAA453345からなる遺伝子産物を含む細胞内物質を検出することを特徴とする乳癌の重篤度検定用の細胞内物質測定機材および測定方法並びに乳癌の重篤度検定方法。

【請求項56】 さらにジーンバンクアクセッション番号がAA406242、H73335、W84753、N71160、AA054669、N32820及びR05667からなる群から選択される少なくとも一つの遺伝子産物を含む細胞内物質を検出することを特徴とする請求項55に記載の乳癌の重篤度検定用の細胞内物質測定機材および測定方法並びに乳癌の重篤度検定方法。

【請求項57】 実質的にジーンバンクアクセッション番号がW84753、H08581、AA045730及びAI250654からなる遺伝子産物を含む細胞内物質を検出することを特徴とする乳癌の再発性検定用の細胞内物質測定機材および測定方法並びに乳癌の再発性検定方法。

【請求項58】 さらにジーンバンクアクセッション番号がAA448641、R78516、R05934、AA629838及びH53037からなる群から選択される少なくとも一つの遺伝子産物を含む細胞内物質を検出することを特徴とする請求項57に記載の乳癌の再発性検定用の細胞内物質測定機材および測定方法並びに乳癌の再発性検定方法。

【請求項59】 実質的にジーンバンクアクセッション番号がAA434397、T83209、N53427、N29639、AA485739、AA425861、H84871、T64312、T59518及びAA037488からなる遺伝子産物を含む細胞内物質を検出することを特徴とする乳癌の再発性検定用の細胞内物質測定機材および測定方法並びに乳癌の再発性検定方

法。

【請求項60】 さらにジーンバンクアクセッション番号がAA406231の遺伝子産物を含む細胞内物質を検出することを特徴とする請求項59に記載の乳癌の再発性検定用の細胞内物質測定機材および測定方法並びに乳癌の再発性検定方法。

【請求項61】 実質的にジーンバンクアクセッション番号がH11482、T64312及びAA045340からなる遺伝子産物を含む細胞内物質を検出することを特徴とする乳癌の再発性検定用の細胞内物質測定機材および測定方法並びに乳癌の再発性検定方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、生体の状態と遺伝子発現の量および／または細胞内物質の量との多変量解析処理並びにそれを基に可能となる測定機材、検定方法などに関するものである。

【0002】

【従来の技術】

2000年6月のヒトゲノムの解読宣言以降、ゲノムに書かれた遺伝情報がどのように発現して機能しているかのを解明するポストゲノム時代に突入したと言われている。ヒトゲノム計画の進展の中で、ゲノム発現状態を測定する方法論も進展してきた。トランスクリプトーム(mRNA)測定手段としてオリゴヌクレオチドアレイやマイクロチップが知られている。またプロテオーム(蛋白質)測定手段として、以前からある2次元電気泳動に加えて、最近では質量分析の方法が進歩してきた。また抗体チップなどの先進の技術も注目されている。これらの測定技術は、生体の状態パラメータを短時間に一挙に測定できることがそれまでの技術と比較して画期的であるといえる。

【0003】

遺伝子発現状態を効率的に測定する技術として次のものがあげられる。トランスクリプトーム(mRNAの総体)を特定するものとして、基盤に複数種のDN

Aを担持し、それに相補的なmRNAを検出するDNAチップが知られている。代表的なDNAチップには、遺伝子チップやDNAマイクロアレイがある。また、プロテオーム（蛋白質の総体）を特定するものには、2次元電気泳動、抗体チップ、質量スペクトルを用いるものがある。またメタボローム（代謝中間体を含めた代謝産物の総体）を測定する手法も質量分析などによって試みられており、進展が見られる。

【0004】

生体内の細胞の状態は遺伝子産物の発現によってよく記述されるため、従来の診断マーカーでは情報が不足している場面でも、精度のより高い診断が可能になるという期待も出てきている。たとえば、次のような研究があげられる。

【0005】

P. O. Brownらは、DNAチップによってリンパ腫患者の細胞のトランスクリプトームを測定し、クラスター解析によって悪性と良性のリンパ腫（DLBCL）を別クラスターに分離した（Nature 403(3), 503-11 (2000)）。しかし、これは因果関係（相関関係）のモデルを得る方法ではなく、どの遺伝子がどの程度重要かを判断できない。

【0006】

A. Alaiyaらは、2次元電気泳動によって子宮がん患者40人の細胞のプロテオームを測定し、うち22人のデータから部分最小自乗法診断モデルを構築し、悪性度を説明した（Int. J. Cancer, 86, 731-36 (2000); Electrophoresis, 21, 1210-17 (2000); 国際公開 WO 00/70340）。その際、全変数モデルにおいて1553変数からloadingの大きな170変数に限定することによって交差検証成績がよくなり（ $Q^2=0.84$ ）、残り18患者の深刻度（3段階）を11/18の比率で正答した。交差検証法がモデル構築の際の指標になるという考えが表明されている。しかし、この方法では、loadingを得る際にまず全変数モデルが成立しなければならない。また、それ以外の変数選択手法が考案されていない。

【0007】

J. Khanらは、DNAチップによって小児がん患者の細胞を測定し、ニューラルネットワークによって悪性度を説明した（Nature Medicine, 7(6), 673-79 (2

001))。小児がん (SRBCT) 患者 88 人のトランスクリプトーム (6567 遺伝子) を測定し、うち 63 人のデータから主成分分析によって 10 次元に圧縮し、次に、人工ニューラルネットワーク診断モデルを構築した。ここで、影響力のある上位遺伝子を交差検証法によって絞り込み、96 遺伝子で最良の成績 (100%) を得た。このモデルで残り 25 人を予測し、93~100% の結果を得た。しかし、この方法でも、影響力を得る際にまず全変数モデルが成立しなければならない。またそれ以外の変数選択手法が考案されていない。10 次元のような少ない変数の場合を扱えるが、変数の数が膨大な場合には適用できない。

【0008】

【特許文献 1】

国際公開公報 WO 00/70340

【非特許文献 1】

P.O. Brownら、Nature 403(3), 503-11 (2000)

【非特許文献 2】

A. Alaiyaら、Int. J. Cancer, 86, 731-36 (2000)

【非特許文献 3】

A. Alaiyaら、Electrophoresis, 21, 1210-17 (2000)

【非特許文献 4】

J. Khanら、Nature Medicine, 7(6), 673-79 (2001)

【非特許文献 5】

P. Gramaticsら、Chemosphere 38(5), 1371-78 (1999)

【0009】

【発明が解決しようとする課題】

従来の診断マーカーでは情報が不足している場面でも、遺伝子発現情報を活用することで、より精度 (解像度) の高い診断が可能になるという期待も出てきている。遺伝子発現状態の測定結果は、膨大な情報量が得られることが従来にはなかった特徴であり、逆に情報量が多いために、効果的なデータ処理なくしてデータの活用はありえない。したがって、有用な知識を獲得するためには効果的な情報処理が欠かせない。前に説明したように、現状ではクラスター解析を中心とす

る方法が用いられているが、主成分分析などの方法も採用されている。クラスター解析や主成分分析は、教師付学習方法ではないため、病状の因果関係（相関関係）のモデルを得ることはできない。すなわち、どの遺伝子がどの程度重要かを解析結果から得ることができないのが難点である。一方、部分最小自乗法は次元圧縮とモデルフィットを同時に行なう強力な多変量解析手法であるが、変数の数が膨大になった場合にしばしば有意な結果が得られない事態に直面する。したがって、膨大な遺伝子発現情報などから有用な知識を獲得できるような効果的な情報処理が望まれている。また、そのような情報処理の結果を基にした効率的な測定機材、検定処理などが期待されている。

【0010】

この発明の目的は、多変量の遺伝子発現情報、細胞内物質情報の効果的な情報処理を提供することである。

また、この発明の目的は、効率的な検定処理を提供することである。

【0011】

【課題を解決するための手段】

本発明に係るデータ解析装置は、生体の状態または時間とともに確率的に発生する生体の状態の変化を目的変数とし、複数の遺伝子発現の量および／または細胞内物質の量を説明変数とする相関モデルを決定するデータ解析装置であって、生体の状態或いはそれを導出するデータまたは時間とともに確率的に発生する生体の状態の変化に関するデータと、複数の遺伝子発現の量および／または細胞内物質の量からなるサンプルの集合を入力する入力手段と、(1)説明変数を選択する選択手段と、(2)部分最小自乗法を実行して交差検証成績を計算する計算手段または前記生体の状態の変化に関するデータにカプラン・マイヤー法又はカトラー・エデラー法による生命表を適用して変化の発生しなかったものの確率を計算して得られた確率を、仮定した分布に基づいた変換または仮定を前提としない変換をし、該変換結果を目的変数とする部分最小自乗法を実行して交差検証成績を計算する計算手段と、(3)前記(2)の計算手段の結果を評価し、説明変数の採用、不採用を判定する評価判定手段とを有し、(4)前記(1)の選択手段と前記(2)の計算手段と前記(3)の評価判定手段とを実行して部分最小自乗法モデルの

少なくとも交差検証成績を独立変数として持つ関数を改善し続けて部分最小自乗法モデルを決定する決定手段とからなる。前記選択手段は、たとえば、説明変数を逐次取捨選択したり、遺伝的アルゴリズムを用いて説明変数を選択する。計算手段は、たとえば、1個のサンプルを逐次除外したり、複数のサンプルを逐次除外して部分最小自乗法を実行して交差検証成績を計算する。評価判定手段は、たとえば、計算手段の結果から、各計算において除外したサンプルの遺伝子発現から予測される生体の状態を示す目的変数値と、前記除外したサンプルの生体の状態を示す目的変数値との誤差の代表値を求め、当該誤差の代表値が小さくなった場合に、その交差検証成績が改善されたと判定し、説明変数を取捨選択しながら交差検証成績の評価判定を繰り返す。あるいは交差検証成績ではなく、少なくとも部分最小自乗法モデルの交差検証成績を独立変数として持つ関数が改善するかどうかを評価判定の基準として用いることもできる。決定手段は、たとえば、選択手段と計算手段と評価判定手段とを繰り返し実行して部分最小自乗法モデルの交差検証成績を改善し続けて部分最小自乗法モデルを決定する。また、選択手段と計算手段とを複数のコンピュータで実行させることもできる。こうして、関連モデルを構成するとき、交差検証成績を基準に最適化させることにより説明変数を取捨選択し、説明変数の数を減らす次元圧縮をして良好なモデルを得る。

【0012】

上述の、仮定した分布に基づいた変換または仮定を前提としない変換は、生体の状態の変化の確率が説明変数の多項式で解析できるようにするために行なうものである。分布を仮定した場合には、確率を対数変換後に負の数にしたものを状態の変化を観測した時間で割るという変換、確率を対数変換後に負の数にしたものをさらに対数にしたものを状態の変化を観測した時間で割るという変換、または確率を1より減じたものをプロビット変換したものを計算しを状態の変化を観測した時間で割るという変換などが考えられる。一方、分布を仮定しない場合にはロジット変換といった方法が考えられる。変換の方法は分布にどのような仮定が成り立つかどうかあるいはなりたたないかどうかを判断することにより、それぞれの場合に応じて適切に選ぶことができる。少なくとも部分最小自乗法モデルの交差検証成績を独立変数として持つ関数としては、たとえば、前記誤差の代

表値と選抜された説明変数の数の関数が考えられ、あるいはその他の独立変数を含むものであってよい。望ましくは、関数は誤差の代表値の単調減少関数であり、説明変数の数の単調減少関数である。計算量を増やさないためには簡単に計算できる関数が望ましい。具体的には $-\text{PRESS} \times \alpha^{NP}$ という関数が考えられる。ここでPRESSは予測残差自乗和であり、NPは採用された説明変数の数であり、 α は1または1より大きい実数である。また、 $-\text{PRESS} \times (NP + \beta)^{\gamma}$ や $-\text{PRESS} \times (\beta - NP)^{-\gamma}$ なる関数も考えられる。ここで、 γ は正の実数である。

【0013】

説明変数の個数を少なくすると、通常の統計的手法または多変量解析手法が適用可能になる。本発明では、部分最小自乗法を用いて得られた説明変数を統計的手法又は多変量解析手法の説明変数として、より良好なモデルを得る。統計的手法又は多変量解析手法としては、重回帰分析法、線型判別分析法、適応最小自乗法、ロジスティック回帰分析法、比例ハザード解析法、マハラビノス距離を用いる判別分析法、kNN法、人工ニューラルネットワークなどが挙げられる。

【0014】

本発明者等は、また、 Q^2 やPRESS値などの交差検証成績に加えて、説明変数の個数を第2の独立変数として含む関数を最適化することで選抜される説明変数を任意に絞り込むことができることを新たに見出した。通常の統計的手法や多変量解析手法では、抽出される説明変数の個数NPの望ましい範囲がサンプル数との兼ね合いで決まっている場合がある。そのような場合、関数を、目的とする選抜数によって任意に変更できる。関数形をたとえば $-\text{PRESS} \times \alpha^{NP}$ とした場合、説明変数の個数を数個から数十個に絞り込むためには通常は定数 α として1.0～3.0の値が望ましい。より望ましくは、 α は1.0～2.0の値となる。他の関数形 $f(\text{PRESS}, NP)$ であっても、実際に選択される説明変数の数MPおよびその時のPRESS値PRESS_MPの周辺で、 $f(\text{PRESS_MP} \div \alpha, MP+1) \approx f(\text{PRESS_MP}, MP)$ となるような関数は、変数選択という点では同様の効果を持つ場合がある。こうして、適当な関数形を用いることにより、望ましい範囲の個数NPの説明変数を選抜できる。このようにして、交差検証成績を用いて決定されたモデルに採

用されている説明変数をさらに絞り込むと、統計的手法又は多変量解析手法によるモデルを構築できる。したがって、その性質が十分解明されている統計的手法又は多変量解析手法を採用して解析を加えることができる。

【0015】

また、目的変数として、時間とともに確率的に発生する生体の状態の変化から導出された量を用いて、時間とともに確率的に発生する生体の状態の変化と複数の遺伝子発現の量および／または細胞内物質の量との相関モデルを決定できる。

「時間とともに確率的に発生する生体の状態の変化」とはたとえば生存時間である。ここで、前述の部分最小自乗法に、カプラン・マイヤー法又はカトラー・エデラー法と、ロジット(logit)変換とを組み合わせる。部分最小自乗法での目的変数は、時間とともに確率的に発生する生体の状態の変化に関するデータにカプラン・マイヤー法又はカトラー・エデラー法による生命表を適用して変化の発生しなかったものの確率を計算し、これをロジット変換した値である。ロジット(logit)値とは、分類分けされたデータの、ある分類の割合(確率) P を基に、次式 $\text{logit} = \log \{P/(1-P)\}$ にて計算される値である。ロジット値を目的変数とする部分最小自乗法を実行して交差検証成績を計算する。こうして、先に説明したのと同様に、部分最小自乗法の交差検証成績を考慮した説明変数の抽出を行って、生存時間解析を行える。

【0016】

説明変数の個数を少なくすると、通常の統計的手法または多変量解析手法が適用可能になる。そこで、決定されたモデルに採用されている説明変数を用い、時間とともに確率的に発生する生体の状態の変化を説明する統計的手法又は多変量解析手法によるモデルを構築する。たとえば、ロジット値を目的変数として求めた説明変数を用いて、他の統計的手法又は多変量解析手法(たとえば比例ハザード法や、パラメトリックな分布にあてはめた回帰分析法)を行なうことによって、より良好なモデルを得ることができる。比例ハザード法とは、Coxによって考案された方法であり、生存率の解析に時間を考慮し、かつ、多変量を扱える。比例ハザード法では、観測されている個々ごとにハザード値と呼ばれる生存率を左右する値があり、それを導く関数がある(モデルが仮定されている)として解析

される。カプラン-マイヤー法は、集団全体または群ごとの生存率の推移を示す。また、パラメトリックな分布とは、ガウスが提案した正規分布から計算された確率分布のことであり、生存時間解析では指数分布、ワイブル分析、対数正規分布が用いられる。指数分布などへの当て嵌めで、数式中に多項式があり、前述の部分最小自乗法の交差検証成績を考慮した説明変数の抽出が適用される。

【 0 0 1 7 】

入力手段で説明変数として入力される複数の遺伝子の発現量および／または細胞内物質の量とは、必ずしも物質の絶対的な濃度の測定値に限定されるものではなく、加工計算された値、相対的な値、間接的に物質量を表す量などでもよい。たとえば、質量スペクトルで蛋白質の発現量を測定することができることを応用して、生体の状態を表わす目的変数と、質量スペクトルとを直接関係づける相関モデルを構築することができる。またAffymetrix社タイプのDNAチップ(ジーンチップ)では、単一のスポットが単一の遺伝子発現を特定するとは限らず、複数個のスポットが集まってはじめて単一の遺伝子発現を特定することもある。ここでもまた、各スポットの測定量を説明変数として、直接、生体の状態を説明する相関モデルを得ることができる。更には、タンパク質の電気泳動パターンの各ピークは単一のタンパク質に帰属できず、複数個のタンパク質の重ねあわせであることも多い。このような場合にも生体の状態を説明する説明変数として各ピーク強度を用いることができる。このことは、上述のAlaiyaらは子宮癌の診断の説明変数として電気泳動パターンのピーク強度を採用していることから明らかである。前述のようにポストシーケンス時代のトランスクリプトーム解析、プロテオーム解析、メタボローム解析という研究分野では、生体(細胞)内の物質を総体として把握することから出発することを特徴とする実験的アプローチが注目されている。ひとつひとつの物質の絶対的定量は必須事項ではなく、これらの実験方法によって定量される物質の量を直接、間接に表現する測定値やその加工計算値が、生体の状態を説明する説明変数と成り得る。また以上の物質量を表現する説明変数以外に、場合によっては問診データなどの他の説明変数を追加すると、さらに有効な解析結果が得られる場合もある。

【 0 0 1 8 】

本発明に係るデータ解析方法は、生体の状態と複数の遺伝子発現の量および／または細胞内物質の量との相関モデルを決定するデータ解析方法であって、生体の状態を目的変数とし、複数の遺伝子発現の量および／または細胞内物質の量を説明変数とするサンプルの集合を入力する入力ステップと、(1)説明変数を選択する選択ステップと、(2)部分最小自乗法を実行して交差検証成績を計算する計算ステップと、(3)前記(2)の計算ステップの結果を評価し、説明変数の採用、不採用を判定する評価判定ステップとを有し、(4)前記(1)の選択ステップと前記(2)の計算ステップと前記(3)の評価判定ステップとを実行して部分最小自乗法モデルの交差検証成績を改善し続けて部分最小自乗法モデルを決定する決定ステップとからなる。前記選択ステップは、たとえば、説明変数を逐次取捨選択したり、遺伝的アルゴリズムを用いて説明変数を選択する。計算ステップは、たとえば、1個のサンプルを逐次除外したり、複数のサンプルを逐次除外して部分最小自乗法を実行して交差検証成績を計算する。評価判定ステップは、たとえば、計算ステップの結果から、各計算において除外したサンプルの遺伝子発現から予測される生体の状態を示す目的変数値と、前記除外したサンプルの生体の状態を示す目的変数値との誤差の代表値を求め、当該誤差の代表値が小さくなった場合に、その交差検証成績が改善されたと判定し、説明変数を取捨選択しながら交差検証成績の評価判定を繰り返す。決定ステップは、たとえば、選択ステップと計算ステップと評価判定ステップとを繰り返し実行して部分最小自乗法モデルの交差検証成績を改善し続けて部分最小自乗法モデルを決定する。また、選択ステップと計算ステップとを複数のコンピュータで実行させることもできる。

【0019】

本発明に係るデータ解析プログラムは、生体の状態と複数の遺伝子発現の量および／または細胞内物質の量との相関モデルを決定する、コンピュータにより実行されるデータ解析プログラムであって、生体の状態を目的変数とし、複数の遺伝子発現の量および／または細胞内物質の量を説明変数とするサンプルの集合を入力する入力ステップと、(1)説明変数を選択する選択ステップと、(2)部分最小自乗法を実行して交差検証成績を計算する計算ステップと、(3)前記(2)の計算ステップの結果を評価し、説明変数の採用、不採用を判定する評価判定ステッ

ブとを有し、(4)前記(1)の選択ステップと前記(2)の計算ステップと前記(3)の評価判定ステップとを実行して少なくとも部分最小自乗法モデルの交差検証成績を独立変数として持つ関数を改善し続けて部分最小自乗法モデルを決定する決定ステップとからなる。

【0020】

上記のデータ解析プログラムにおいて、前記選択ステップは、たとえば、説明変数を逐次取捨選択したり、遺伝的アルゴリズムを用いて説明変数を選択する。計算ステップは、たとえば、1個のサンプルを逐次除外したり、複数のサンプルを逐次除外して部分最小自乗法を実行して交差検証成績を計算する。評価判定ステップは、たとえば、計算ステップの結果から、各計算において除外したサンプルの遺伝子発現から予測される生体の状態を示す目的変数値と、前記除外したサンプルの生体の状態を示す目的変数値との誤差の代表値を求め、少なくとも当該誤差の代表値を独立変数として持つ関数である当該誤差の誤差の代表値の単調減少関数の値が小さくなった場合に、その交差検証成績が改善されたと判定し、説明変数を取捨選択しながら交差検証成績の評価判定を繰り返す。決定ステップは、たとえば、選択ステップと計算ステップと評価判定ステップとを繰り返し実行して少なくとも部分最小自乗法モデルの交差検証成績を独立変数として持つ関数を改善し続けて部分最小自乗法モデルを決定する。また、選択ステップと計算ステップとを複数のコンピュータで実行させることもできる。さらには、上記の説明変数の選択において、たとえば、初期状態では説明変数を全く含まないか、或いは、初期状態では全説明変数を含むこともできる。

【0021】

上記のデータ解析プログラムにおいて、上記の生体の状態は、たとえば病気のタイプあらかず測定値、病気の重篤度をあらかず測定値、病気のタイプをあらかず医療診断の結果、病気の重篤度をあらかず医療診断の結果、あるいはそれらを2次加工した数値である。例えば後の実施例で示すように、患者の生存時間を予測することは、QOLを含めた治療計画や人生設計などを判断する上で重要な情報をもたらすものであり、社会的に価値のある診断モデルを提供することができる。また癌の再発可能性を予測することは、QOLを考慮した治療計画を立案し、医

師または当の患者が選択の判断をするうえで、貴重な情報をもたらすものである。

【 0 0 2 2 】

また、本発明は、前記で決定された相関モデル及び予測対象のサンプルについて当該モデルにおいて採用された説明変数を入力する入力手段と、入力された該説明変数に基づいて該サンプルの生体の状態を予測判定する予測判定手段からなるデータ解析装置、前記で決定された相関モデル及び予測対象のサンプルについて当該モデルにおいて採用された説明変数を入力する入力ステップと、入力された該説明変数に基づいて該サンプルの生体の状態を予測判定する予測判定ステップからなるデータ解析方法及び前記で決定された相関モデル及び予測対象のサンプルについて当該モデルにおいて採用された説明変数を入力する入力ステップと、入力された該説明変数に基づいて該サンプルの生体の状態を予測判定する予測判定ステップからなるデータ解析プログラムも包含する。

【 0 0 2 3 】

本発明に係るコンピュータにより読取可能な記録媒体は、上記のいずれかのプログラムを記録する。

【 0 0 2 4 】

本発明に係るびまん性大細胞型Bリンパ腫の重篤度検定用の細胞内物質測定機材および測定方法並びにびまん性大細胞型Bリンパ腫の重篤度検定方法は、実質的にジーンバンクアクセッション番号がU15085、M23452、X52479、U70426、H57330及びS69790からなる遺伝子群の発現を検出する。さらに、ジーンバンクアクセッション番号がU03398、M65066、AK001546、BC003536、X00437、U12979、H96306、AA830781及びAA804793からなる群から選択される少なくとも一つの遺伝子の発現を検出してもよい。

【 0 0 2 5 】

また、本発明に係る乳癌の重篤度検定用の細胞内物質測定機材および測定方法並びに乳癌の重篤度検定方法は、実質的にジーンバンクアクセッション番号がAA598572、AA703058及びAA453345からなる遺伝子産物を含む細胞内物質を検出する。さらに、ジーンバンクアクセッション番号がAA406242、H73335、W84753、N

71160、AA054669、N32820及びR05667からなる群から選択される少なくとも一つの遺伝子産物を含む細胞内物質を検出してもよい。

【0026】

また、本発明に係る乳癌の再発性検定用の細胞内物質測定機材および測定方法並びに乳癌の再発性検定方法は、実質的にジーンバンクアクセッション番号がW84753、H08581、AA045730及びA1250654からなる遺伝子産物を含む細胞内物質を検出する。さらに、ジーンバンクアクセッション番号がAA448641、R78516、R05934、AA629838及びH53037からなる群から選択される少なくとも一つの遺伝子産物を含む細胞内物質を検出してもよい。

【0027】

また、本発明に係る乳癌の再発性検定用の細胞内物質測定機材および測定方法並びに乳癌の再発性検定方法は、実質的にジーンバンクアクセッション番号がAA434397、T83209、N53427、N29639、AA485739、AA425861、H84871、T64312、T59518及びAA037488からなる遺伝子産物を含む細胞内物質を検出する。さらに、ジーンバンクアクセッション番号がAA406231の遺伝子産物を含む細胞内物質を検出してもよい。

【0028】

また、本発明に係る乳癌の再発性検定用の細胞内物質測定機材および測定方法並びに乳癌の再発性検定方法は、実質的にジーンバンクアクセッション番号がH11482、T64312及びAA045340からなる遺伝子産物を含む細胞内物質を検出する。

【0029】

細胞内物質測定機材としては、DNAマイクロアレイ、ジーンチップ、オリゴDNA型のDNAチップ、電気化学DNAチップ(ECAチップ)、繊維型DNAチップ、磁性ビーズDNAチップ(PSS)、糸巻きDNAチップ(PSS)、などのDNAチップ、マイクロアレイ、抗体チップ、測定用試薬キットなどが挙げられる。また、上記の機材を適宜組み込んだ測定機械であってもよい。

【0030】

【発明の実施の形態】

以下、添付の図面を参照して本発明の実施の形態を説明する。

以下に、選択された生体の状態と遺伝子発現の量および／または細胞内物質の量との相関モデルの決定について説明する。ここで、遺伝子発現の用語は、mRNA発現(トランスクリプトーム)や、mRNAによる翻訳の結果として生じる蛋白質(プロテオーム)を含むものとして用いる。また、細胞内物質の量とはここではたとえば、代謝中間体を含めた代謝産物全部であるメタボロームを意味する。たとえば、トランスクリプトーム(mRNA)やプロテオーム(蛋白質)の解析において、各サンプルデータは、生体の状態と遺伝子発現の量などからなる。各サンプルはたとえば1000個以上の膨大な遺伝子発現の量を含む。生体の状態は、たとえば病気のタイプまたは病気の診断指標であるが、より一般的には生体情報であればよい。「病気の診断指標」には、病気の進行度合いのほか、病気のタイプ、重篤度、深刻度などの表現で表わされるものも含む。ここで、遺伝子発現の量などの測定データは膨大な情報量からなるので、コンピュータを用いた効率的な多変量解析が必要である。

【0031】

データ収集において、予めいくつかのサンプルについて生体の状態(たとえば診断指標)を判定し、また、そのサンプルされたものから細胞液を獲得し、その細胞液中の多くの遺伝子産物の発現の量などを測定する。本発明の実施の形態のデータ解析では、こうして得られた遺伝子産物の発現の量などと生体の状態(たとえば診断指標)を入力し、相関モデル(たとえば部分最小自乗法モデル)を得る。ここで、コンピュータによる多変量解析プログラムを用いて、診断指標を目的変数とし、遺伝子発現の量および／または細胞内物質の量を説明変数とする因果関係型の解析を行なって、各説明変数の重要性や影響度に関する情報を得る。また、前記目的変数は、必ずしも測定値そのものである必要はなく、ロジット変換を行なった値や群を表す離散値を用いても良く、その場合、より有意な解析結果を得ることもできる。

【0032】

本発明者らは、遺伝子発現による医療診断という分野において、データ解析における交差検証(cross validation)の成績を少なくとも独立変数のひとつとして持つ関数を最適化するように変数を選択することによって良好な相関モデル(

たとえば部分最小自乗法モデル) が得られることを見出した。交差検証法では、手持ちのデータを複数群に分割し、その一部のデータ群(訓練集合)だけを使ってフィットしたモデルを用いて残る別のデータ群(テスト集合)を予測することによって、モデルの予測力を試す。通常の部分最小自乗法(PLS)においては潜在変数の次元選択に交差検証法が用いられているが、ここでは、部分最小自乗法において、潜在変数を1次元に固定し、1以上の入力変数(説明変数)を逐次取捨選択しながら、交差検証成績(たとえば平方和の予測誤差)を少なくとも独立変数のひとつとして持つ関数を最適化した。その結果、全変数を採用した場合には有意な相関モデルを得られなかった場合にも、良好でかつ予測力のある相関モデルが得られることが判明したのである。この交差検証法を用いた変数選択の逐次取捨選択により、安定な相関モデルが得られる。また本発明者らは、関数形を適切に設定することによって説明変数を絞り込むことにより、部分最小自乗法以外の統計学又は多変量解析の良好な相関モデルを得ることが可能となり、それぞれ生体の状態を記述する目的変数にふさわしい相関モデルを得ることができることを見出した。なお、ここでいう「最適化」とは、交差検証成績が、説明変数を取捨選択するための、そのときの解析条件の範囲で、改善がみられなくなるまで改良したことを意味しており、交差検証成績がすべての説明変数の組合せの中で最適なものを見出したという意味ではない。この変数選択手法を用いると、病状を決定する因子を少数に特定し、廉価な診断用材料(DNAチップ、抗体チップ、DNA含有ベクターなど)を設計でき、それ自体独自の価値を持つものである。また、この変数選択手法は、予め設定される各種の変数選択条件と共に運用することが可能である。

【0033】

上に述べたように、説明変数は、交差検証成績を基準に逐次取捨選択される。ここで、取捨選択のため、交差検証成績を少なくとも独立変数のひとつとして持つ関数を用いる。説明変数を追加する場合は、その説明変数について、前記関数が改善されなかったと判定された場合には当該説明変数を除外し、改善されたと判定された場合には当該説明変数を追加する。また、説明変数を除外する場合は、その説明変数について、前記関数が改善されなかったと判定された場合には当

該説明変数を除外せず、改善されたと判定された場合には当該説明変数を除外する。ここで、1以上の説明変数を選択した場合に、交差検証成績評価は次のように進める。n個のサンプルからいくつかのサンプルを逐次除外して部分最小自乗法モデルを求め、各モデルにおいて除外したサンプルの遺伝子発現の量から予測される生体の状態を示す目的変数と、除外したサンプルの生体の状態を示す目的変数との各々の誤差の代表値を求める。「代表値」とは、和、平均、最大値、中位値、最頻値などのデータを特徴づける値をいう。そして、当該誤差の代表値を少なくともひとつの独立変数とする関数が小さくなった場合に、交差検証成績が改善されたと判定し、当該説明変数を追加または削除する。この交差検証成績評価を、説明変数を取捨選択しながら逐次繰り返して、前記関数を改善し続ける。改善されなくなれば交差検証成績を最適化したとして説明変数の取捨選択を終了する。その結果、取捨選択により絞り込んだ数の説明変数からなる最適な部分最小自乗法モデルが得られる。具体的には、計算手段において計算される交差検証成績の数値指標として予想残差自乗和(PRESS)を採用し、評価判定手段において予想残差自乗和の値が説明変数あたり一定の閾値以下の比率で小さくなる場合に、その説明変数を採用すると判定することにより、上記の処理は実行可能である。

【0034】

因果関係型の解析手法においてはオーバーフィット (over fitting) を避けるための工夫が必要となる。ここでいうオーバーフィットとは、説明変数が多すぎるためにたまたま予測結果と実績とが一致するものの、本当の相関関係をとらえ損なっているため、モデルフィットに用いたデータ以外に予測能力を持たないことをいう。ここでは、相関モデルとして部分最小自乗法を用いるが、部分最小自乗法は次元圧縮とモデルフィットを同時に行なう強力な多変量解析手法であり、オーバーフィットの問題に比較的強いとされている。しかし遺伝子発現状態解析のように膨大な変数を扱う場合には、有意な結果が得られない事態に直面する。従来技術として説明したAlaiyaやKhanの手法は全変数モデルが有意に成立することを前提としているので、変数の絞り込みには一般的には適用できない。これに対し、本発明では、交差検証予測結果を最適にするように変数を絞り込むことによ

り、オーバーフィットを減らすことができた。また、本発明は、前記Khanの手法とは異なり、主成分分析などの前処理を介さない方法である。従来技術では、説明変数が膨大な場合には、有意なモデルを得ることができないことから、予め、全説明変数を基にたとえば、主成分分析などで次元圧縮する前処理をし、これによって得られた説明変数によって解析する方法が用いられる。しかし、この方法では、構成したモデルで予測を行なうためには、モデル構成の基となった全説明変数が必ず必要となり、たとえば、説明変数が遺伝子発現の量であれば、診断用遺伝子チップに担持する遺伝子としては、モデル構成に用いた遺伝子の全てが必要となるか、または別の手法を用いて変数選択することが必要となる。一方、本発明においては、説明変数の選択によって説明変数を絞り込んでいるので、たとえば、説明変数が遺伝子発現の量であれば、診断用遺伝子チップに担持する遺伝子は、選択された説明変数に相当する遺伝子を担持すれば良いことになる。

【0035】

なお、Todeschiniらは、有機化合物の大気中の分解を予測するため、遺伝的アルゴリズムによって交差検証成績を最適化するように変数選択を行ない、重回帰モデルを得ている (P. Gramatica, V. Consonni & R. Todeschini, *Chemosphere* 38(5), 1371-78 (1999))。53化合物と175記述子でモデル構築を行ない ($Q^2=0.79$)、7変数が選択され、98化合物の予測を行なった ($Q^2=0.75$)。交差検証成績を最適化するように変数選択を行なっている点では、本実施形態と同様の手法である。しかし、重回帰モデルを採用しているために、説明変数の選択過程を通じて選択される変数は少数個にとどまらざるを得ず、複数の遺伝子発現の量および／または細胞内物質の量の解析には適用できない。本発明者らの調査した範囲では、 Q^2 やPRESS値を最適化する方法では、選抜される説明変数は百程度から数百程度にわたり、重回帰モデルでは解析が不能となる。またTodeschiniらは、説明変数を絞り込むための有効な方法について言及していない。これは、もともとの説明変数の候補がたかだか175個であり、説明変数を絞り込むために特別の工夫をする必要がないからである。遺伝子発現解析の分野はこれとは全く異なり、数十から数百のサンプル数に対して、数百から数千、数万の説明変数候補が存在する。したがってこれまでとは異なる工夫が必要となる。

【0036】

本実施形態では、生体の状態と複数の遺伝子発現の量および／または細胞内物質の量との相関モデルを決定するとき、交差検証成績を少なくとも独立変数のひとつとして持つ関数を最適化させるように説明変数を逐次追加・除外することによって、説明変数を選抜して、良好な相関モデルを得る。このようなアプローチの優位性は、下記の実施例から推測されるように、次のとおりである。

- 1) 病気や生体现象の背後で働いている重要な遺伝子やメカニズムを推定／特定でき、理解が深まる。
- 2) 重要な遺伝子産物や細胞内物質だけに絞った廉価な診断用材料（DNAチップ、抗体チップなど）の設計が可能になる。

【0037】

本実施形態では、交差検証成績を少なくとも独立変数のひとつとして持つ関数を最適化するように説明変数を段階的に取捨選択するが、たとえば具体的には、ステップワイズ(step wise)法に代表される説明変数を選択する選択手段と、リーブ・ワン・アウト(leave-one-out)法に代表される交差検証法に部分最小自乗法を適用して計算する計算手段と、前記計算手段の結果を評価し、説明変数の採用、不採用を判定する評価判定手段とを組合せて用いる。すなわち、 m 個の説明変数の中から1以上の説明変数を選択し、次いで、部分最小自乗法を実行して交差検証成績を計算し、さらに、該計算結果を評価して、選択した説明変数の採用、不採用を判定する。この評価判定では、計算手段の結果から、各計算において除外したサンプルの遺伝子発現から予測される生体の状態を示す目的変数値と、前記除外したサンプルの生体の状態を示す目的変数値との誤差の代表値を求め、少なくとも当該誤差の代表値を独立変数として持つ関数である当該誤差の誤差の代表値の単調減少関数の値が小さくなった場合に説明変数の取捨選択を判定する。このように、選択手段と計算手段と評価判定手段とを用いて、少なくとも部分最小自乗法モデルの交差検証成績を独立変数として持つ関数を改善し続けて、その改善がみられなくなるまで改良し、部分最小自乗法モデルを決定する。なお、本実施形態では、サンプルを1個づつ逐次除外している(リーブ・ワン・アウト法)が、その代わりに、複数のサンプルを除外して交差検証成績を評価してもよ

い(リーブ・n・アウト法)し、また、Khan et al.により用いられた3分割法(three-fold)等の他の方法を用いることもできる。3分割法では、説明変数をランダムにシャッフルして3つのグループに分ける。その中の2つのグループを用いてモデルを構成し、残りの1つのグループでモデルを評価する。また、説明変数の選択方法としてはステップワイズ法、非線形アルゴリズム(たとえば遺伝的アルゴリズムなど)を用いてもよく、変数選択に関して予め何らかの条件が分っていれば、それに応じて探索範囲を限定できる。

【0038】

次に、データの収集と解析について具体的に説明する。図1は、遺伝子発現解析システムを示す。データ収集のため、予めいくつかのサンプルについて診断指標(たとえば病気のタイプないし進行度合いを含む)を判定し、また、そのサンプルされたものから細胞液を獲得し、DNAチップを用いてその細胞液中の多くの遺伝子産物の発現の量を測定する。測定には、共焦点型レーザスキャナ(たとえばAffymetrix社、428アレイスキャナ)10を用いる。吸光度によりmRNAの量が測定される。このデータ収集は公知の方法である。測定データは、コンピュータ12に送られ解析される。コンピュータ12は、CPU14を備えた通常の構成のコンピュータであり、それに接続される記憶装置(たとえばハードディスク装置)16の記録媒体(たとえばハードディスク)には、測定データ18や解析ソフト20が格納される。この解析ソフト20を用いてデータ18が解析され、生体の状態と遺伝子発現の量などとの相関モデルが決定される。

【0039】

なお、説明変数の選択と、交差検証法に部分最小自乗法を適用する計算とを複数のコンピュータで実行させてもよい。交差検証予測の計算を複数個のコンピュータに分散させることで計算を加速することができる。

【0040】

図2は、コンピュータ12により実行される、生体の状態と遺伝子発現の量などとの相関モデルを得るためのデータ解析ソフト20のフローチャートを示す。ここでは簡単に説明するため、少なくとも部分最小自乗法モデルの交差検証成績を独立変数として持つ関数として-PRESSを採用しているが、発明の範囲を限定す

るものでなく、実施例 2～5 においては別の関数を採用している。まず、相関モデル作成用のデータを入力する (S 1 0)。データはたとえば DNA チップを用いて収集したものである。入力データ (サンプル集合) は、それぞれ目的変数 (たとえば診断指標) と m 個 (たとえば 2 0 0 0 個) の説明変数 (たとえば遺伝子発現の量) からなる。また、場合によっては、上述のデータ (訓練集合) 以外に、テスト集合のデータを入力する。ここでテスト集合とは交差検証の評価のためのデータ群を意味するのではなく、モデル決定が終了した後にモデルの予測力をテストするためのデータ群である。

【 0 0 4 1 】

まず、初期設定として、選択された説明変数の数を 0 とし、交差検証成績 CV の最良値 CV_0 を $-\infty$ とする (S 1 2)。次に、説明変数の選択を行う。まず、説明変数を指す番号 i を 1 とし (S 1 4)、第 i 変数 (遺伝子発現の量) を仮に採用して (S 1 6)、部分最小自乗法を実行し、交差検証成績 CV を計算する (S 1 8、図 3 参照)。ここで、リーブ・ワン・アウト処理を用いる。これは、たとえば 5 0 個のサンプルからなる訓練集合において、1 番から 5 0 番の全てを順次 1 個づつ除いて残りの 4 9 個のサンプルで予測した結果と、その時除いた 1 個の結果とを比較し、その誤差が大きい場合に、仮に選択した説明変数 (第 i 変数) が適していないと判断する手法である。もし、得られた成績 CV が現在の最良値 CV_0 より最適化されれば (S 2 0 で YES)、第 i 変数を採用し、かつ、成績 CV を新しい最良値 CV_0 に更新する (S 2 2)。しかし、得られた成績 CV が最良値 CV_0 より大きくなければ (S 2 0 で NO)、第 i 変数を採用しない (S 2 4)。そして、ステップ S 1 4 に戻り、同様の処理を繰り返す。この処理を交差検証成績 CV が改善されなくなる (S 2 6 で NO) まで繰り返す。ここで、相関モデルに採用する説明変数については 1 つずつ段階的に増加 (追加) または減少 (除外) して成績 CV を評価判定している。すなわち、全体としての合致度合いがよくなるように各説明変数を解析に加えるかどうかを逐次判定しながら、説明変数の取捨選択を行い、これを、全体としての合致度合いがよくなるまで繰り返す。以上の処理で改善があると、ふたたびステップ S 1 4 の初め ($i=1$) に戻り、それまでに選択されている説明変数を基に、さらに説明変数の選択を繰り返す。なお、ここで

はモデルの予測力を判断するために、訓練集合とテスト集合とに予め分割しておいたデータ集合を用いてデータ解析しており、上述の解析は、訓練集合を用いて行なった結果であるので、この結果からテスト集合について予測を行い、実測データとの合致度を評価（S28）している。このような評価は必ずしも必要でないが、予測力を判断するには有効である。

【0042】

図3は、リーブ・ワン・アウト処理を含む交差検証成績CVの計算（図2、S18）のフローチャートを示す。ここで、選択された変数について交差検証成績が計算される。まず、PRESSの初期値を0とする（S180）。次に、 n 個の集合内のサンプルを指す番号 j を1とし（S182）、第 j サンプル以外の $n-1$ 個のサンプルで部分最小自乗法を実行し（S184）、第 j サンプルの目的変数を予測する（S186）。差の自乗を計算してPRESSに加算する（S190）。次に番号 j を1増加し（S182）、同様の処理をおこなう。これを番号 $j=n$ まで各サンプルについて繰り返す。得られたPRESSは、1個のサンプルを順次除外して計算した予測値と実測値との差の平方和であり、予測誤差を表わす量である。この予測残差自乗和PRESSの符号を変えたものを交差検証成績CVとする（S192）。

【0043】

本実施形態では、交差検証法を用いて、入力変数（説明変数）を段階的に1つづつ追加・除外しながら、交差検証成績（ $CV=-PRESS$ ）を最適化する。ここで、説明変数の段階的な追加・除外の内容を理解しやすくするため、以下で、さらに具体的に5つのモデル構築手法について説明する。これらは、説明変数の逐次的な選択の手順が異なる。

【0044】

図4は、第1のモデル構築手法を示す。データ集合においてどの説明変数も選択されていない状態を初期状態とする（S112）。次に、1番目の説明変数から最後（ m 番目）の説明変数までの未だ選択されていない説明変数ごとに逐次、その説明変数を選択した場合に交差検証成績が改善するかどうかを、リーブ・ワン・アウト処理を用いた交差検証成績評価ステップ（S118）を繰り返しながら

判定(S 1 2 0)し、改善する場合にはその説明変数を追加する(S 1 1 4 ~ S 1 2 4)。そのような改善と追加がなくなる(S 1 2 6でNO)まで、1番目の説明変数から上記逐次判定操作を繰り返す。

【0 0 4 5】

さらに詳しく説明すると、まず、初期設定として、選択された説明変数の数 N_P を0とし、交差検証成績CVの最良値 CV_0 を $-\infty$ とする(S 1 1 2)。次に、説明変数の選択を行う。まず、変数 i を1とし(S 1 1 4)、第 i 変数を仮に採用する(S 1 1 6)。ただし、第 i 変数がすでに採用されていれば(S 1 1 5でYES)、ステップS 1 1 4に戻る。次に、部分最小自乗法を実行し、交差検証成績CVを計算する(S 1 1 8)。ここで、リーブ・ワン・アウト処理を用いる。もし、得られた成績CVが現在の最良値 CV_0 より最適化されれば(S 1 2 0でYES)、第 i 変数を採用し、かつ、成績CVを新しい最良値 CV_0 に更新する(S 1 2 2)。しかし、得られた成績CVが最良値 CV_0 より大きくなければ(S 1 2 0でNO)、第 i 変数を採用しない(S 1 2 4)。そして、ステップS 1 1 4に戻り、同様の処理を繰り返す。この処理を交差検証成績CVが改善されなくなる(S 1 2 6でNO)まで繰り返す。以上の処理で改善があると、ふたたびステップS 1 1 4に戻り、新しいループを開始する。ここで、それまでに選択されている変数を基に、さらに変数の選択を繰り返す。こうして、データ集合を用いて選択された変数を用いた相関モデルが得られる。

【0 0 4 6】

図5は、第2のモデル構築手法を示す。この手法では、全ての説明変数が選択されている状態を初期状態とする(S 2 1 2)。次に、1番目の説明変数から最後(m 番目)の説明変数までの選択されている説明変数ごとに逐次、その説明変数を除外した場合に交差検証成績が改善するかどうかを、リーブ・ワン・アウト処理を用いた交差検証成績評価ステップ(S 2 1 8)を繰り返しながら判定(S 2 2 0)し、改善する場合にはその説明変数を除外する(S 2 1 4 ~ S 2 2 4)。そのような改善と除外がなくなる(S 2 2 6でNO)まで、1番目の説明変数から上記逐次判定操作を繰り返す。

【0 0 4 7】

さらに詳しく説明すると、まず、初期設定として、選択された説明変数の数 N_P を m とし、交差検証成績 CV の最良値 CV_0 を $-\infty$ とする (S 2 1 2)。すなわち、すべての説明変数を選択する。次に、説明変数の選択を行う。まず、変数 i を 1 とし (S 2 1 4)、第 i 変数を仮に除外する (S 2 1 6)。ただし、第 i 変数がすでに除外されていれば (S 2 1 5 で YES)、ステップ S 2 1 4 に戻る。部分最小自乗法を実行し、交差検証成績 CV を計算する (S 2 1 8)。ここで、リーブ・ワン・アウト処理を用いる。もし、得られた成績 CV が現在の最良値 CV_0 より最適化されれば (S 2 2 0 で YES)、第 i 変数を除外し、かつ、成績 CV を新しい最良値 CV_0 に更新する (S 2 2 2)。しかし、得られた成績 CV が最良値 CV_0 より大きくなければ (S 2 2 0 で NO)、第 i 変数を除外しない (S 2 2 4)。そして、ステップ S 2 1 4 に戻り、同様の処理を繰り返す。この処理を交差検証成績 CV が改善されなくなる (S 2 2 6 で NO) まで繰り返す。以上の処理で改善があると、ふたたびステップ S 2 1 4 に戻り、新しいループを開始する。ここで、それまでに選択されている変数を基に、さらに変数の選択を繰り返す。こうして、データ集合を用いて選択された変数を用いた相関モデルが得られる。

【 0 0 4 8 】

図 6 は、第 3 のモデル構成手法を示す。この手法は、第 1 と第 2 の手法の直列的な組合せである。まず、どの説明変数も選択されていない状態を初期状態とする (S 1 1 2)。次に、1 番目の説明変数から最後 (m 番目) の説明変数までの未だ選択されていない説明変数ごとに逐次、その説明変数を選択した場合に交差検証成績が改善するかどうかを、リーブ・ワン・アウト処理を用いた交差検証成績評価ステップを繰り返しながら判定し、改善する場合にはその説明変数を追加選択し、そのような改善と追加がなくなるまで 1 番目の説明変数から上記逐次判定操作を繰り返す (S 1 1 4 ~ S 1 2 6)。次に、1 番目の説明変数から最後 (m 番目) の説明変数までの選択されている説明変数ごとに逐次、その説明変数を除外した場合に交差検証成績が改善するかどうかを、リーブ・ワン・アウト処理を用いた交差検証成績評価ステップを繰り返しながら判定し、改善する場合にはその説明変数を除外し、そのような改善と除外がなくなるまで、1 番目の説明変数から上記逐次判定操作を繰り返す (S 2 1 4 ~ S 2 2 6)。

【 0 0 4 9 】

図 7 は、第 4 のモデル構築手法を示す。この手法は、第 3 の手法の変形である。まず、どの説明変数も選択されていない状態を初期状態とする (S 1 1 2)。次に、1 番目の説明変数から最後 (m 番目) の説明変数までの未だ選択されていない説明変数ごとに逐次、その説明変数を選択した場合に交差検証成績が改善するかどうかを、リーブ・ワン・アウト処理を用いた交差検証成績評価ステップ (S 1 1 8) を繰り返しながら判定 (S 1 2 0) し、改善する場合にはその説明変数を追加選択する (S 1 1 4 ~ S 1 2 4)。そのような改善と追加がなくなる (S 1 2 6 で N O) まで、1 番目の説明変数から上記逐次判定操作を繰り返す。次に、1 番目の説明変数から最後 (m 番目) の説明変数までの選択されている説明変数ごとに逐次、その説明変数を除外した場合に交差検証成績が改善するかどうかを、リーブ・ワン・アウト処理を用いた交差検証成績評価ステップ (S 2 1 8) を繰り返しながら判定 (S 2 2 0) し、改善する場合にはその説明変数を除外する (S 2 1 4 ~ S 2 2 4)。そのような改善と除外がなくなる (S 2 2 6 で N O) まで、1 番目の説明変数から上記逐次判定操作を繰り返す。上記逐次判定追加改善ステップまたは上記逐次判定除外改善ステップで少なくとも一度改善があれば (S 2 2 7 で Y E S)、ステップ S 1 1 2 に戻り、上記操作 (S 1 1 2 ~ S 2 2 7) を繰り返す。これを改善がなくなる (S 2 2 7 で N O) までおこなう。

【 0 0 5 0 】

図 8 は、第 5 のモデル構築手法を示す。この手法は、第 1 と第 2 のスキームの並列的な組合せである。どの説明変数も選択されていない状態を初期状態とする (S 1 1 2)。次に、1 番目の説明変数から最後 (m 番目) の説明変数までの説明変数ごとに逐次、その説明変数が選択されていない場合にはその説明変数を選択した場合に交差検証成績が改善するかどうかを、リーブ・ワン・アウト処理を用いた交差検証成績評価ステップ (S 1 1 8) を繰り返しながら判定 (S 1 2 0) し、改善する場合にはその説明変数を追加する (S 1 1 4 ~ S 1 2 4)。また、選択する説明変数ごとに、その説明変数がすでに選択されている場合には、その説明変数を除外した場合に交差検証成績が改善するかどうかを、リーブ・ワン・アウト処理を用いた交差検証成績評価ステップ (S 2 1 8) を繰り返しながら

判定 (S 2 2 0) し、改善する場合にはその説明変数を除外する (S 2 1 6 ~ S 2 2 4)。そのような改善と追加または除外がなくなる (S 1 2 6 で NO) まで、1 番目の説明変数から上記逐次判定操作を繰り返す。

【 0 0 5 1 】

次に、第 4 のモデル構築手法 (図 7) を適用した場合を、表 1 のデータ集合を例として説明する。このデータ集合に対して、部分最小自乗法による解析を用いて相関モデルを求める。表 1 のデータでは、サンプルの数 n は 1 0 であり、また、説明を容易にするため、説明変数の数 m は 1 9 と少なくしている。表 1 において、 p_1 は目的変数を表わし、 $p_2 \sim p_{20}$ は説明変数を表わす。(ただし表 1 では、表示の便宜のため、 p_{16} 以降のデータを省略している。) 第 4 手法 (図 7) のステップ S 1 1 4、S 2 1 4 とは異なり、説明変数を表わす i は p_{20} から p_2 まで逆に逐次処理することとした。CV 評価値としてここでは予測残差自乗和 (PRESS) を採用した。PRESS が小さいほど、CV 評価値はよい。初期状態では、採用された説明変数の数 NP は 0 であり、 $PRESS = \infty$ ($CV_0 = -\infty$) である。

【 0 0 5 2 】

【表 1】

表 1 1 0 個のサンプルのデータ

#	p_1	p_2	p_3	p_4	p_5	p_6	p_7	p_8	p_9	p_{10}	p_{11}	p_{12}	p_{13}	p_{14}	p_{15}
1	0.713	0.105	0.782	0.425	0.164	0.023	0.696	0.543	0.333	0.691	0.336	0.668	0.017	0.061	0.5
2	0.133	0.009	0.071	0.002	0.793	0.872	0.092	0.391	0.63	0.241	0.517	0.389	0.166	0.841	0.1
3	0.545	0.193	0.765	0.334	0.109	0.538	0.578	0.652	0.38	0.501	0.729	0.91	0.865	0.389	0.8
4	0.752	0.915	0.472	0.999	0.798	0.363	0.622	0.487	0.353	0.967	0.778	0.484	0.517	0.982	0.0
5	0.9	0.407	0.534	0.816	0.806	0.42	0.572	0.957	0.12	0.696	0.833	0.051	0.377	0.849	0.4
6	0.455	0.587	0.721	0.53	0.252	0.434	0.882	0.486	0.741	0.243	0.893	0.947	0.462	0.952	0.2
7	0.427	0.652	0.515	0.426	0.764	0.592	0.595	0.595	0.551	0.606	0.416	0.163	0.316	0.718	0.6
8	0.042	0.902	0.274	0.899	0.402	0.469	0.688	0.945	0.746	0.912	0.97	0.515	0.368	0.514	0.4
9	0.935	0.276	0.936	0.101	0.54	0.356	0.899	0.71	0.924	0.792	0.486	0.329	0.501	0.076	0.5
10	0.54	0.021	0.505	0.224	0.724	0.431	0.071	0.968	0.482	0.322	0.773	0.543	0.353	0.107	0.9

【 0 0 5 3 】

【表2】

表2 表1のデータについての10の段階での変数選択結果

0		∞	-
1	追加	p20 0.111	p20
2	追加	p18 0.090	p18 & p20
3	追加	p16 0.073	p16 & p18 & p20
4	追加	p10 0.073	p10 & p16 & p18 & p20
5	追加	p6 0.062	p6 & p10 & p16 & p18 & p20
6	追加	p3 0.060	p3 & p6 & p10 & p16 & p18 & p20
7	追加	p12 0.055	p3 & p6 & p10 & p12 & p16 & p18 & p20
8	除外	p20 0.053	p3 & p6 & p10 & p12 & p16 &
9	除外	p10 0.050	p3 & p6 & p12 & p16 & p18
10	追加	p13 0.048	p3 & p6 & p12 & p13 & p16 & p15

【0054】

先に述べたように、変数はp20からp2まで逆の順で処理する。表2は、表1のサンプルについて、左端の数字は、変数の取捨選択で改善がみられた10の段階を示す。なお、0は初期状態を意味する。次の列の「追加」と「除外」は、追加のループと除外のループの処理であることを意味する。次の列の変数は、追加または除外された変数を示す。次の列は、交差検証成績(PRESSをサンプル数で割ったもの)を示す。右端の列は、その段階で選択されている変数を示す。

【0055】

初期状態では、変数は全くない状態であり、PRESSは ∞ である。表2に示すように、最初、p20を説明変数として採用すると、PRESS=0.111となり、初期値に比べて改善されるので、説明変数p20の追加を実施する。次に、変数p19を加えてp19とp20の2つを説明変数とすると、PRESS=0.129となり改善をもたらさないで、p19は追加しない。次に、説明変数p18を加えるとPRESS=0.090となり、改善するので、p18を追加し、p18とp20を説明変数とする。以下同様に表2に示すように続く。(ここで、p10を追加採用するのは、小数点以下4桁目で改善されているた

めである。)説明変数 $p_{20} \sim p_2$ の1回目のループを終了した時点で、説明変数が p_3 、 p_6 、 p_{10} 、 p_{16} 、 p_{18} および p_{20} となり、 $PRESS=0.60$ となる。2回目のループでは、説明変数 p_{12} が追加され、 $PRESS=0.55$ となる。3回目のループでは追加による改善がなく、ひとまず $S114 \sim S126$ の追加処理を終了し、 $S214$ に移る。この時点での部分最小自乗法のフィットならびにリーブ・ワン・アウト予測状況は表3のとおりである。

【0056】

表3は、10のサンプルについて、表2の7で示す段階まで処理が進んだ時点での部分最小自乗法のフィットならびにリーブ・ワン・アウト予測状況を示す。ここで、モデル予測とリーブ・ワン・アウト予測のそれぞれにおいて、計算値と実測値との誤差を示す。さらに、その下側に、誤差の自乗平均、相関係数 R の自乗および予測相関係数 Q の自乗を示す。

【0057】

【表 3】

表 3 表 2 の段階 7 での処理結果

#	実測値	モデル予測値		リーフワンアウト予測	
		計算値	誤差	計算値	誤差
1	0.713	0.757	-0.044	0.693	0.020
2	0.133	-0.056	0.189	-0.051	0.184
3	0.545	0.497	0.048	0.480	0.065
4	0.752	0.646	0.106	0.495	0.257
5	0.900	0.687	0.214	0.557	0.343
6	0.455	0.489	-0.034	0.512	-0.057
7	0.427	0.624	-0.198	0.672	-0.245
8	0.042	0.349	-0.307	0.517	-0.475
9	0.935	0.865	0.070	0.782	0.153
10	0.154	0.197	-0.044	0.285	-0.132
		0.093	0.024	0.055	
		$R^2=0.744$	$Q^2=0.407$		

【0058】

次に、S 2 1 4 から始まる除外処理の 1 回目のループにおいて、説明変数 p10 と p20 を除外することが改善をもたらした。2 回目のループでは改善がなく、S 2 1 4 ～ S 2 2 6 を終了するが、S 2 2 7 の判断により再度 S 1 1 2 に戻る。次に、追加処理の 1 回目のループにおいて、p13 の追加だけが改善をもたらしたが、続く除外処理の 1 回目のループでは、改善がなかった。もう一度 S 1 1 2 に戻り、ステップ S 1 1 4 ～ S 1 2 6 およびステップ S 2 1 4 ～ S 2 2 6 では改善がなくなったのを確認して、処理を終了した。こうして選択された説明変数は、p3、p6、p12、p13、p16 および p18 の 5 個であり、PRESS=0.048 となった。詳細は表 4 のとおりである。

【0059】

表4は、表2の段階10まで処理が進んだ時点での部分最小自乗法のフィットならびにリーブ・ワン・アウト予測状況を示す。

【0060】

【表4】

表4 表2の段階10での処理結果

#	実測値	モデル予測		リーブワンアウト予測	
		計算値	誤差	計算値	誤差
1	0.713	0.771	-0.058	0.663	0.050
2	0.133	-0.013	0.146	0.041	0.092
3	0.545	0.610	-0.065	0.595	-0.050
4	0.752	0.524	0.228	0.380	0.372
5	0.900	0.696	0.205	0.543	0.357
6	0.455	0.591	-0.137	0.623	-0.168
7	0.427	0.638	-0.211	0.696	-0.269
8	0.042	0.189	-0.147	0.268	-0.226
9	0.935	0.841	0.094	0.756	0.179
10	0.154	0.209	-0.055	0.294	-0.140
		0.093	0.022	0.048	
		$R^2=0.765$		$Q^2=0.482$	

【0061】

なお、説明変数の数が多い時に強いとされる部分最小自乗法であるが、p20～p2の全てを説明変数として採用した場合には、表5に示すように、PRESS=0.124となった。すなわち、リーブ・ワン・アウト処理は、平均値からの誤差(0.093)よりも悪い成績をもたらす。

【0062】

【表 5】

表 5 全ての説明変数を採用した場合の処理結果

#	実測値	モデル予測		リーブワンアウト予測	
		計算値	誤差	計算値	誤差
1	0.713	0.712	0.001	0.527	0.186
2	0.133	-0.073	0.206	0.222	-0.090
3	0.545	0.561	-0.016	0.538	0.007
4	0.752	0.656	0.096	0.351	0.402
5	0.900	0.691	0.209	0.432	0.469
6	0.455	0.519	-0.064	0.562	-0.107
7	0.427	0.583	-0.156	0.629	-0.203
8	0.042	0.430	-0.388	0.724	-0.682
9	0.935	0.794	0.140	0.480	0.454
10	0.154	0.182	-0.029	0.457	-0.303
<hr/>					
	0.093		0.029		0.124
	$R^2=0.684$			$Q^2=-0.330$	

【 0 0 6 3 】

【実施例】

次に、実施例を挙げて本発明をさらに詳細に説明するが、本発明はこれらの例によって何ら限定されるものではない。

【 0 0 6 4 】

実施例 1： 部分最小自乗法の交差検証成績を考慮した特徴抽出による DLBCL 患者のデータ解析

P. O. Brownらのホームページ (<http://llmpp.nih.gov/lymphoma/>) より入手した 28 名の DLBCL (リンパ腫) 患者のデータを、20 名のデータからなる訓練集合と 8 名のデータからなるテスト集合に分けた。目的変数に生存月数を採用し、説明変数には 18432 スポットのうち、28 データにおいて ch1、ch2 と

もに正の数となる12832スポットの $\log(\text{ch1}/\text{ch2})$ 値を採用した。

【0065】

訓練集合において部分最小自乗法 (PLS) のモデル決定を試みた。12832変数全てを用いて部分最小自乗法の解析をしたところ、リーブ・ワン・アウト予測は有意($Q^2 > 0.5$)にはならなかった。次にリーブ・ワン・アウト予測誤差が最小になるように説明変数を段階的に1つつ増減した。モデル構成手法としては前述の第3のモデル構成手法において説明変数の追加及び除外の順番並びにリーブ・ワン・アウト処理におけるサンプルの除外の順番が異なるほかは同様な方法を用いた。すなわち、どの説明変数も選択されていない状態を初期状態とする(S112)。次に、最後(m番目)の説明変数から最初(1番目)の説明変数までの未だ選択されていない説明変数ごとに逐次、その説明変数を選択した場合に交差検証成績が改善するかどうかを、リーブ・ワン・アウト処理(ここでは、最後(n番目)のサンプルから最初(1番目)のサンプルを逐次除外した)を用いた交差検証成績評価ステップを繰り返しながら判定し、改善する場合にはその説明変数を追加選択し、そのような改善と追加がなくなるまでm番目の説明変数から上記逐次判定操作を繰り返す(S114~S126)。次に、最後(m番目)の説明変数から最初(1番目)の説明変数までの選択されている説明変数ごとに逐次、その説明変数を除外した場合に交差検証成績が改善するかどうかを、リーブ・ワン・アウト処理(ここでも最後(n番目)のサンプルから逐次除外した)を用いた交差検証成績評価ステップを繰り返しながら判定し、改善する場合にはその説明変数を除外し、そのような改善と除外がなくなるまで、最後(m番目)の説明変数から上記逐次判定操作を繰り返す(S214~S226)。その結果、有意なモデル($R^2=0.988$ 、 $Q^2=0.895$ 、 $NP=342$)を得た。図9は、このデータについての最小自乗法成績を示す。図9において、ひし形は訓練集合のデータ(20人)を示し、三角は、それらについての交差検証成績のデータを示す。また、四角はテスト集合のデータ(8人)を示す。得られた部分最小自乗法モデルは、テスト集合のうち、4/8をきわめて良好に、また1/8を良好に予測するものであった。

【0066】

なお、上述の多変量解析によるデータ解析では、扱ったサンプルはDNAチップを用いて得たデータであった。しかし、このデータ解析は、DNAチップを用いて得たデータに限定されるものではなく、蛋白質発現量、細胞内物質の量などのデータに対しても有用であろうことは容易に推測されることである。

【 0 0 6 7 】

以下の実施例 2～7 では、部分最小自乗法を用いて選抜した少ない個数の説明変数について、通常の統計的手法または多変量解析手法（比例ハザード法、重回帰分析、適応最小自乗法、ロジスティック回帰分析法、線型判別分析法など）を適用する。

【 0 0 6 8 】

実施例 2： 部分最小自乗法の交差検証成績を考慮した特徴抽出と比例ハザード解析による 240 名の DLBCL 患者の生存時間解析

RosenwaldらがWeb上 (<http://llmpp.nih.gov/DLBCL/>) で公開している 240 名の DLBCL (びまん性大細胞型 B リンパ腫) のデータセットをダウンロードして用いた。全データを訓練集合として利用した。スポットパターンで x_1 または x_2 が 0 となるものを除いた 7399 スポットについて $\log(x_1/x_2)$ を計算して説明変数とした。本実施例では実施例 1 と異なり、生存時間として観測打ち切り時間と死亡時間とが混在していることを考慮してカプラン・マイヤー (Kaplan-Meier) 法による生命表を適用して事象発生時点での生存確率 (P_{KM}) を求め、ロジット変換 ($\log(P_{KM}/1-P_{KM})$) した値を目的変数とした。カプラン・マイヤー法による生命表は集団としての生存確率を示すが、ここでは、個人 j を含む集団としての事象発生時点での残存確率 (変化の発生しなかったものが残存する確率) を個人 j の事象発生時点での残存時間に読み代えるという新規な考え方をを用いている。また、この確率をロジット変換して、変化の発生傾向を表現するロジット値に変換して、目的変数とした。訓練集合内の交差検証はリーブ・ワン・アウト法によって行ない、 $PRESS \times 1.02^{N_P}$ が小さくなるようにパラメータを逐次取捨選択して部分最小自乗法モデルを得た。ここで、交差検証成績 ($CV = -PRESS$) の代わりに、少なくとも交差検証成績を独立変数として持つ関数の 1 つである関数 $-PRESS \times 1.02^{N_P}$ を改善して部分最小自乗法モデルを得た。ここで $PRESS$ はリーブ・

ワン・アウト予測の残差自乗和であり、NPは、選択された説明変数の数である。

【 0 0 6 9 】

図7のフロー中の交差検証成績CVを $-\text{PRESS} \times 1.02^{NP}$ と読み換えて、処理を実行することにより、下記の19個の遺伝子の発現が説明変数として選抜された。ここでdata IDはWebデータ元でのID番号を示す。またACCESSIONはGenBankのアクセシオン番号であり、アクセシオン番号の無い行はデータ元でのみ明らかとなっている遺伝子(Unknown)ないしESTであり、論文記載の方法によって入手することができる。

【 0 0 7 0 】

ACCESSION	data ID	comment
U03398	#(27876)	tumor necrosis factor (ligand) superfamily, member 9
M65066	#(27394)	protein kinase, cAMP-dependent, regulatory, type I, beta
--	#(27104)	(Unknown)
AK001546	#(25048)	Homo sapiens cDNA FLJ10684 fis, clone NT2RP3000220
--	#(31372)	(Unknown)
U15085	#(28178)	major histocompatibility complex, class II, DM beta
BC003536	#(24983)	hypothetical protein MGC10796
--	#(16113)	(Unknown)
M23452	#(16822)	small inducible cytokine A3
	#(24433)	(Unknown)
X00437	#(27480)	T cell receptor beta locus
U12979	#(24377)	activated RNA polymerase II transcription cofactor 4
X52479	#(17773)	protein kinase C, alpha

H96306	#(16578)	bone marrow stromal cell antigen 1
U70426	#(19255)	regulator of G-protein signaling 16
AA830781	#(33358)	EST
AA804793	#(25022)	EST
H57330	#(26383)	EST
S69790	#(27184)	WAS protein family, member 3

【 0 0 7 1 】

これらの遺伝子の発現を説明変数の候補として比例ハザード(hazard)解析を試みた。比例ハザード法とは、生存率の解析に時間を考慮した統計的手法である。解析の実行はプログラムパッケージ JMP (JMP Sales SAS Campus Drive Cary, N C 27513 USA)を用いて行なった。変数削除基準として $P \geq 0.05$ を採用した変数減少法によって更に絞り込んだ結果、14 遺伝子の発現からなる以下の比例ハザード式が得られた。ここでGenbank (ジーンバンク) のアクセシオン番号ないし data IDで示される各項は、各遺伝子の $\log(x_1/x_2)$ 値であり、またPは統計的な有意性が成り立たない危険率である。この式の右辺から求められるハザード(hazard)値が大きいほど、死亡傾向が大きい。

【 0 0 7 2 】

$$\begin{aligned} \text{hazard} = & 0.370 \#(27104) + 0.589 \text{ AK001546} - 0.366 \#(31372) - 0.276 \text{ U15085} \\ & - 0.307 \#(16113) + 0.409 \text{ M23452} - 0.350 \#(24433) - 0.297 \text{ X00437} \\ & + 0.321 \text{ U12979} - 0.585 \text{ X52479} - 0.457 \text{ U70426} + 0.561 \text{ AA830781} \\ & - 0.430 \text{ H57330} + 0.433 \text{ S69790} \end{aligned}$$

$$P < 0.0001$$

【 0 0 7 3 】

Rosenwaldらは、単相関の比例ハザード解析を行なって、5群(17 遺伝子)の診断指標を選抜している。図10に、本実施例で得られたハザード値(Hazard)とRosenwaldらの診断指標がどの程度、生存時間を説明できているかを比較した。Rosenwaldらの5群のパラメータを同時に用いた比例ハザード式ではProliferationパラメータが $P > 0.05$ で統計的に有意でないなどの問題を有していたため、これを除く4群のパラメータを同時に含めたハザード値も比較のために掲載した。こ

こで、菱形は死亡した人または打ち切った人のデータを示し、四角は生存している人のデータを示す。

【 0 0 7 4 】

これらの診断指標のうち、本実施例で求めたハザード値と生存時間との相関は際立って明白である。即ちハザード値は生存時間につれて減衰しており、大きなハザード値の患者は長く生きることが出来ないことが示されている。一方、Rose nwaldらの指標はいずれも生存時間を診断するには不十分なものである。数百、数千という数のパラメータの中から効率的に最適のパラメータセットを見出すことは比例ハザード解析だけではできないことである。しかし以上のようにカプラン・マイヤー法、ロジット変換、部分最小自乗法の交差検証成績を考慮した特徴抽出、比例ハザード解析を組み合わせることで、従来に無い、有効な診断指標を得ることができた。統計学的に異質なモデルをこのように組み合わせることによってこのような良好な結果が得られたことは意外でもあり、興味深いことであった。患者の生存時間を予測することは、QOLを含めた治療計画や人生設計などを判断する上で重要な情報をもたらすものであり、本実施例で求められた診断モデルは社会的に価値のあるものである。

【 0 0 7 5 】

また、変数削除基準として $P \geq 0.001$ を採用した変数減少法によって更に絞り込むと、7遺伝子の発現からなる以下の比例ハザード式が得られた。このように、変数削除基準を変えることにより、選択される説明変数の数を制御できる。

【 0 0 7 6 】

$$\begin{aligned} \text{hazard} = & -0.426 \text{ U15085} + 0.350 \text{ M23452} - 0.521 \text{ X52479} \\ & - 0.450 \text{ U70426} - 0.586 \text{ H57330} + 0.476 \text{ S69790} \end{aligned}$$

【 0 0 7 7 】

図 1 1 は、右辺を計算して求められるハザード値を縦軸とし、生存時間を横軸としたプロットを示す。図 1 0 と同様に、図 1 1 において、菱形は死亡した人または打ち切った人のデータを示し、四角は生存している人のデータを示す。

【 0 0 7 8 】

実施例 3 : 部分最小自乗法の交差検証成績を考慮した特徴抽出と比例ハザー

ド解析による40名の乳癌患者の生存時間解析

SorleらがWeb上(http://genome-www.stanford.edu/breast_cacer/mopo_clinical/)で公開している乳癌患者のデータセットをダウンロードして用いた。全データを訓練集合として利用した。データセットの大部分は、タイプA, Bという2種類のDNAチップで測定されたそれぞれ40名、24名の患者よりなるが、ここではタイプAのデータを用いた。生存時間データより実施例2と同様にロジット値を求め、目的変数とした。説明変数としては、データに欠測のある遺伝子を除いた6891件のLOG_RAT2N_MEAN値を採用した。そして、少なくとも交差検証成績を独立変数として持つ関数の1つである、交差検証成績と説明変数NPの関数 $\text{PRESS} \times 1.13^{NP}$ が小さくなるようにパラメータを逐次取捨選択して部分最小乗法モデルを得た。図7のフロー中の交差検証成績CVを $-\text{PRESS} \times 1.13^{NP}$ と読み換えて、処理を実行することにより、下記の10個の遺伝子の発現が説明変数として選抜された。

【0079】

ACCESSION	comment
AA406242	(guanosine monophosphate reductase)
AA598572	(spleen tyrosine kinase)
H73335	(Homo sapiens mRNA full length insert cDNA clone EUROIMAG
E	980547)
W84753	(Homo sapiens cDNA FLJ13510 fis, clone PLACE1005146)
AA703058	(myeloperoxidase)
N71160	(cytochrome c oxidase subunit Vib)
AA453345	(a protein tyrosine kinase)
AA054669	(Homo sapiens, clone IMAGE:3611719, mRNA, partial cds)
N32820	(ESTs, Weakly similar to ALU1_HUMAN ALU SUBFAMILY J SEQUE
NCE	
	CONTAMINATION WARNING ENTRY [H. sapiens])
R05667	(suppressor of potassium transport defect 3)

【0080】

これらを説明変数の候補として、比例ハザード解析において変数削除基準として $P \geq 0.05$ を採用した変数減少法を試み、7遺伝子の発現からなる以下の比例ハザード式が得られた。ここでアクセション番号で示される各項はそれぞれの遺伝子のLOG_RAT2N_MEANである。

【0081】

hazard = -0.821 AA406242 +1.556 AA598572 -1.074 H7335 +1.418 W84753
-1.290 AA703058 +2.182 N71160 +0.828 AA453345

$P < 0.0001$ 変数の $P < 0.05$

【0082】

図12に、右辺を計算して求められるハザード値を縦軸とし、生存時間を横軸としたプロットを示す。ここでもハザード値が優れた診断指標となることが示されている。図12において、菱形は死亡した人または打ち切った人のデータを示し、四角は生存している人のデータを示す。

【0083】

変数削除基準として $P \geq 0.001$ を採用した変数減少法によって更に絞り込んだ。これにより、3遺伝子の発現からなる以下の比例ハザード式が得られた。このように、変数削除基準を変えることにより、説明変数の数を制御できた。

【0084】

hazard = 1.453 AA598572 -1.473 AA703058 +1.071 AA453345

【0085】

図13は、右辺を計算して求められるハザード値を縦軸とし、生存時間を横軸としたプロットを示す。ここで、菱形は死亡した人のデータを示し、四角は生存している人のデータを示す。

【0086】

実施例4： 部分最小自乗法の交差検証成績を考慮した特徴抽出と重回帰分析による40名の乳癌患者の再発予測解析

SorleらのDNAチップAで6891遺伝子の発現が測定された40名の患者をデータセットとして用いた。再発の有無を目的変数として、 $\text{PRESS} \times 1.10^{N_P}$ が小

さくなるようにパラメータを逐次取捨選択して11遺伝子の発現からなる部分最小自乗法モデルを得た。

【0087】

ACCESSION	comment
AA434397	integrin, beta 5
T83209	ESTs
N53427	KIAA1628 protein
N29639	cytidine monophosphate-N-acetylneuraminic acid hydroxylase
AA485739	major histocompatibility complex, class II, DR beta 5
AA425861	enoyl Coenzyme A hydratase 1, peroxisomal
H84871	Ste-20 related kinase
T64312	prostate cancer overexpressed gene 1
T59518	solute carrier family 2, (facilitated glucose transporter) member 8
AA406231	KIAA0381 protein
AA037488	prolactin

【0088】

次に、選抜された遺伝子発現を説明変数とし、再発の有無を目的変数として、通常の多変数解析法の一つである重回帰分析によって判別分析を実行した。解析の実行はプログラムパッケージJMPを用いて行なった。変数削除基準として $P \geq 0.15$ を採用した変数減少法によってさらに絞り込んだ結果、10遺伝子の発現からなる以下の重回帰式が得られた。この式で計算されるOLS値が正の時は再発の可能性が高く、負の時は低い。

【0089】

$$\begin{aligned} \text{OLS} = & -0.215 \text{ AA434397} + 0.227 \text{ T83209} - 0.209 \text{ N53427} + 0.139 \text{ N29639} \\ & + 0.165 \text{ AA485739} + 0.133 \text{ AA425861} - 0.084 \text{ H84871} - 0.193 \text{ T64312} \\ & + 0.237 \text{ T59518} + 0.176 \text{ AA037488} - 0.278 \end{aligned}$$

$R^2=0.84797$ 、 判別正解率 97.5%

【0090】

上式に含まれる各パラメータをそれぞれ1つ用いて判別分析式を作成した場合のP値及び決定係数を以下の表6に示す。

【0091】

【表6】

表6

Accession No.	P value	決定係数(R^2)
AA434397	0.0334	0.090273
T83209	0.0601	0.066005
N53427	0.0004	0.268678
N29639	0.0552	0.069483
AA485739	0.0421	0.080733
AA425861	0.0861	0.05122
H84871	0.087566	0.087566
T64312	0.0004	0.263207
T59518	0.0066	0.157196
AA037488	0.0031	0.187627

単独では有意とはならない($P>0.05$)パラメータが3つ存在し、また、どのパラメータも決定係数が小さい。従って、パラメータを1つずつ吟味するだけでは、上式のような良好な判別式は得られなかった。また数百、数千という数のパラメータの中から効率的に最適のパラメータセットを見出すことは重回帰分析だけではできないことである。しかし、以上のように、部分最小自乗法の交差検証成績を考慮して特徴抽出することにより、従来に無い、有効な診断指標を得ることができた。乳癌の再発可能性を予測することは、QOLを考慮した治療計画を立案し判断するうえで、社会的に求められているところのものである。

【0092】

実施例5： 部分最小自乗法の交差検証成績を考慮した特徴抽出と適応最小自

乗法による40+24名の乳癌患者の再発予測解析

DNAチップのタイプA(40名)とタイプB(24名)に共通する3448遺伝子に限って解析を試みた。PRESS $\times 1.17^{NP}$ が小さくなるようにパラメータを逐次取捨選択して部分最小自乗法モデルを得た。選抜された遺伝子発現を説明変数とし、適応最小自乗法によって判別分析を実行した結果、次式が得られた。次式で計算されるALS値が0.5より大きいと再発の危険性が存在する。

【0093】

$$ALS = 0.31 \text{ H11482} - 0.29 \text{ T64312} - 0.32 \text{ AA045340} + 0.01$$

$$R^2 = 0.65, \text{ eps} = 0.13, \text{ 判別正解率 } 90.0\%$$

【0094】

下記の表7にみるように、H11482は単相関では有意ではなく、他の変数と同時に用いることで初めて把握できたパラメータである。また、表8は、上式を用いてタイプBの患者を予測した結果である。本判別式の感度=81.8%、特異度=53.8%となり、 $\chi^2=3.233$ (5% $P<10\%$)、予測判別正解率=66.7%、という統計的に有意な結果を得た。タイプA、BはDNAチップの構成の相違に基づく測定誤差が存在すると思われるデータであるにもかかわらず、タイプAで訓練したモデルでタイプBの予測に危険率10%以下で成功したことは勇気付けられる結果である。

【0095】

また、PRESS $\times 1.12^{NP}$ が小さくなるように選んだ場合には、以下の遺伝子の発現を説明変数とする部分最小自乗法モデルを得た。

【0096】

$$\text{H11482, T64312, R99749, T65211, AA427625, AA455506}$$

【0097】

これらを説明変数の候補として、リーブ・ワン・アウトを指標にして、さらに絞り込んだ結果、次の判別式を得た。

$$ALS = 0.53 \text{ H11482} - 0.31 \text{ T64312} - 0.33 \text{ R99749} - 0.26 \text{ AA455506} + 0.10$$

$$R^2 = 1.00, \text{ eps} = 0.10, \text{ 判別正解率 } 100.0\%$$

【0098】

パラメータを1つずつ吟味するだけでは、上式のような良好な判別式は得られなかった。また数百、数千という数のパラメータの中から効率的に最適のパラメータセットを見出すことは、適応最小自乗法、ロジスティック回帰分析、その他の判別分析手法だけではできないことである。しかし、以上のように、部分最小自乗法の交差検証成績を考慮して特徴抽出することにより、従来に無い、有効な診断指標を得ることができた。

【0099】

【表7】

表7 パラメータの交絡作用

パラメータ	R	Nmis(140)
H22482	0.861	14
T64312	0.607	8
AA045340	0.572	9
T64312 & AA045340	0.716	6
H11482 & T64312 & AA045340	0.804	4

【0100】

【表8】

表8 タイプBの24患者の予測

観察値	予測値	頻度
-	-	7
+	-	2
-	+	6
+	+	9

【0101】

実施例6

部分最小自乗法の交差検証成績を考慮した特徴抽出とロジスティック回帰分析法または線型判別分析法による40+24名の乳癌患者の再発予測解析

実施例 5 での 1 つめの適応最小自乗法による解析をロジスティック回帰分析法に置き換えた場合、次の判別式が得られた。

【0 1 0 2】

$$\text{LORA} = 7.92 \text{ H11482} - 5.69 \text{ T64312} - 6.41 \text{ AA045340} - 9.73$$

$$R^2 = 0.63, x_2 = 35.00 (P < 0.0001), \text{判別正解率 } 90.0\%$$

【0 1 0 3】

右辺で求められる L O R A 値が正の場合には再発の危険性が存在する。係数の比率や相関係数は実施例 5 の適応最小自乗法の場合と異なるものの、各患者の識別結果は全く同一であった。またタイプ B の患者を予測した結果も表 7 と同じになった。

【0 1 0 4】

次に、実施例 5 での適応最小自乗法による解析を線型判別分析に置き換えて解析して、次の判別式が得られた。

【0 1 0 5】

$$\text{LDA} = 2.45 \text{ H11482} - 2.35 \text{ T64312} - 2.56 \text{ AA045340} - 4.03$$

$$\text{判別正解率 } 80.0\%$$

【0 1 0 6】

右辺で求められる L D A 値が正の場合には再発の危険性が存在する。係数の比率や相関係数は、実施例 5 の適応最小自乗法の場合と異なり、各患者の識別結果も若干異なったが、概ね同一であった。また、タイプ B の患者を予測した結果も表 7 と同じになった。

【0 1 0 7】

以上の実施例 4, 5, 6 では、乳癌の再発の有無を目的変数としている。したがって、部分最小自乗法の交差検証成績を考慮して特徴抽出する方法が、目的変数が名義尺度や順序尺度などのデータである場合にも有効であることが示された。なお、名義尺度とは、対象（サンプル）をある分類に属するかどうかを測り分けるときの分類で、分類の間に大小や順序はない。また、順序尺度とは、対象の特定の分類について測り分けるときの分類であり、分類の間に大小、高低といった順序がある。

【0108】

実施例 7

部分最小自乗法の交差検証成績を考慮した特徴抽出と比例ハザード解析による 40 名の乳癌患者の再発時間解析

実施例 4 と同じデータを用いて、再発の時系列データを基に実施例 2 と同様の方法で求めたロジット値を目的変数として、 $\text{PRESS} \times 1.15^{N/P}$ が小さくなるようにパラメータを逐次取捨選択して 9 遺伝子の発現からなる部分最小自乗法モデルを得た。これらの遺伝子発現の測定値を説明変数として比例ハザード解析において変数削除基準として $P \geq 0.05$ を採用した変数減少法を試み、8 遺伝子からなる、以下の比例ハザード式が得られた。

【0109】

$$\begin{aligned} \text{hazard} = & 1.122 \text{ AA448641} - 1.781 \text{ R78516} - 1.434 \text{ R05934} + 2.165 \text{ W84753} \\ & - 1.923 \text{ AA629838} + 2.665 \text{ H08581} + 1.875 \text{ AA045730} + 1.269 \text{ AI250654} \\ & P < 0.0001 \end{aligned}$$

【0110】

図 1 4 は、右辺を計算して求められるハザード値を縦軸とし、再発時間を横軸としたプロットを示す。ここで、菱形は再発しない人のデータを示し、四角は再発している人のデータを示す。ここでもハザード値が優れた診断指標となっており、生存時間に限らず、時間とともに確率的に発生する生体の状態の変化を解析する手法として、本発明の手法が有効であることが示されている。

【0111】

変数削除基準として $P \geq 0.005$ を採用した変数減少法によって更に絞り込んだ場合には、4 遺伝子の発現からなる以下の比例ハザード式が得られた。

【0112】

$$\text{hazard} = 1.559 \text{ W84753} + 2.265 \text{ H08581} + 1.473 \text{ AA045730} + 1.237 \text{ AI250654}$$

【0113】

図 1 5 は、右辺を計算して求められるハザード値を縦軸とし、再発時間を横軸としたプロットを示す。ここで、菱形は再発しない人のデータを示し、四角は再発している人のデータを示す。

【 0 1 1 4 】

実施例 8 : Genbank アクセション番号 H11482、T64312、AA045340 を含む乳癌再発性診断用 DNA チップの作成と測定

実験医学別冊「ゲノム機能研究プロトコール」(ISBN4-89706-932-7 C3047) p34-38 記載の関直彦、永杉友美、東孝典、吉川勉、鈴木収、村松正明らの方法に順じて DNA チップの作成と測定を行なう。Genbank アクセション番号 H11482、T64312、AA045340 の cDNA を用いる。

【 0 1 1 5 】

プローブ用の各 PCR 産物をエタノール(和光純薬, Cat#057-00456)で沈殿させ、 $2\mu\text{g}/\mu\text{l}$ となるように DDW で調整する。ニトロセルロース(GibcoBRL Cat#41051-012) $4\text{mg}/\text{ml}$ の DMSO 溶液を等量加え、よく混和させて 100°C で 5 分間熱変性を行ない、氷上で急冷する。次いで室温に戻し、DNA スポッター SPBI 02000(日立ソフトエンジニアリング)を用いてカルボジイミドスライドガラス(日清紡)へのスポッティングを速やかに行なう。スポットの乾燥を確認し、Ultraviolet crosslinker(アマシャムファルマシアバイオテック社)を用いて $60\text{mJ}/\text{cm}^2$ で紫外クロスリンク処理を行ない、ガラスラックに立てて室温保存する。

3%BSA、0.2M NaCl、0.1M Tris(PH 7.5)、0.05% Triton X-100 よりなるブロッキング液に上記マイクロアレイを浸け、約 30 分間放置する。次いで、ガラスに付着している溶液をよく切り、 37°C で乾燥させる。TE バッファー(PH 8.0, ニッポンジーン Cat #316-90025)で 3 回軽く洗い、プレートホルダーに入れて軽く遠心(1000 rpm, 1 分間)して余分な水分を除去する。

【 0 1 1 6 】

次に、乳腺正常株 SV-40 及び乳癌細胞株 MCF-7、MDA-MB-468 又は T-47-D の各細胞液より、TRIZOL(GibcoBRL, Cat#15596-018)、Oligotex dT30<Super> (TaKaRa, Cat#W9021A)を用いてマニュアルに従って、mRNA を精製する。mRNA $2\mu\text{g}$ を $6.4\mu\text{l}$ の DEPC 処理し、DDW に溶かし、Oligo dT プライマー $9\mu\text{l}$ 、5X SuperScript II バッファー(GibcoBRL, Cat#18089-011) $6\mu\text{l}$ 、DTT(SuperScript の付属) $3\mu\text{l}$ 、50 X dNTP $0.6\mu\text{l}$ 、Cy3-dUTP(アマシャムファルマシアバイオテック Cat# PA53022)又は Cy5-dUTP (アマシャムファルマシアバイオテック Cat# PA55022) $3\mu\text{l}$

、SuperScript II 2 μ lよりなる溶液を加え、42℃で2時間反応させる。途中1時間経過時点で、SuperScript IIを1 μ lを追加する。1.5 μ lアルカリバッファー(1N NaOH / 20mM EDTA)を加え、65℃で10分間反応させ、TEバッファーを270 μ l、1N HClを1.5 μ l加えて、Cy3, Cy5ラベルの反応液を2つ纏めて1本のMicrocon-YM-30(Millipore/Amicon, Cat#42410)に移す。10,000 rpmで上のカップに残る液量が約10 μ lになるまで遠心を続け、カップを通りぬける液を別のチューブに移し替え、その後、上のカップにTE バッファー500 μ l、Human Cot-1 DNA(GibcoBRL Cat#15279-011) 20 μ gを加え、再び液量が10 μ l以下になるまで遠心を続ける。3,000 rpmで3分間遠心し、蛍光標識したDNAを回収する。DDWとyeast RNA(Sigma, Cat#R7125) 50 μ g、poly(A) (ロッシユダイアグノスティクス, Cat#108 626) 50 μ gを加えて20 μ lにし、PCR用のチューブに移し換え、さらに4.25 μ l 20 X SSC(GibcoBRL, Cat#15553-035)と0.75 μ l 10% SDS(GibcoBRL, Cat#15553-035)を加え、PCR用の機器で100℃、1分間熱変性させ、次いで、室温で30分間放置して、ゆっくり冷却する。

【0117】

蛍光標識したDNAの全量をカバーガラスにのせ、泡が入らないように注意しながら前記マイクロアレイにかぶせ、水で濡らしたキムタオルを底に敷いたハイブリダイゼーションチェンバーに入れて密閉する。毎分2～4サイクルで軽く振とうさせながら、65℃で一晩ハイブリダイズする。ハイブリダイゼーションチェンバーからマイクロアレイを取り出し、カバーガラスがのったままの状態ですぐに2 X SSC / 0.1% SDS溶液の中に入れ、5分間シェイキングし、カバーガラスが自然にはがれるのを待つ。カバーガラスがはがれたところでマイクロアレイをスライドガラスラックに入れ、もう一度2 X SSC / 0.1% SDS溶液中で5分間軽く振とうして洗う。さらに0.2 X SSC / 0.1% SDS 40℃で5分間2回洗い、0.2 X SSCでリンスする。マイクロアレイを別の乾いたプレパレートケースに移し、マイクロタイタープレート用の遠心機で軽く遠心して(1000 rpm, 1分室温)マイクロアレイ上の水分を除く。ScanArray4000(GSI luminonics社)でシグナルを読み込み、解析ソフトにはQuant Array (GSI luminonics社)およびChip Space(日立ソフトウェアエンジニアリング)を用いる。

【0118】

【発明の効果】

生体の状態と複数の遺伝子発現の量および／または細胞内物質の量との相関モデルを決定するとき、説明変数の選択と交差検証法とを用いて変数を絞り込むことができる。これにより、良好でかつ予測力のある多変量解析モデル（相関モデル）が得られる。特に遺伝子発現の量のように、説明変数の数がたとえば1000以上と膨大な場合に有用である。変数の数を少なくすることにより、病気や生体现象の背後で働いている重要な遺伝子やメカニズムを推定／特定でき、理解が深まる。また、重要な遺伝子産物や細胞内物質だけに絞った廉価な診断用材料（DNAチップ、DNA含有ベクター、抗体チップなど）を設計し、提供できる。

【0119】

また、時間とともに確率的に発生する生体の状態の変化から導出された量を目の変数として用いて、時間とともに確率的に発生する生体の状態の変化と複数の遺伝子発現の量および／または細胞内物質の量との相関モデルを決定できる。

【0120】

また、部分最小自乗法を用いて説明変数の個数を少なくすると、通常の統計的手法または多変量解析手法が適用可能になる。

【図面の簡単な説明】

- 【図1】 遺伝子発現解析システムのブロック図
- 【図2】 解析ソフトのフローチャート
- 【図3】 交差検証成績CVの計算のフローチャート
- 【図4】 変数選択の第1モデル構築手法のフローチャート
- 【図5】 変数選択の第2モデル構築手法のフローチャート
- 【図6】 変数選択の第3モデル構築手法のフローチャート
- 【図7】 変数選択の第4モデル構築手法のフローチャート
- 【図8】 変数選択の第5モデル構築手法のフローチャート
- 【図9】 最小自乗法モデルの成績を示すグラフ
- 【図10】 DLBCL患者の生存時間と診断指標のプロット各種比較の図
- 【図11】 実施例2のDLBCL患者の生存時間診断指標のプロット

【図 1 2】 実施例 3 の乳癌患者の生存時間診断指標のプロット

【図 1 3】 実施例 3 の乳癌患者の変数削除基準として $P \geq 0.0005$ を採用したときの生存時間診断指標のプロット

【図 1 4】 実施例 7 の乳癌患者の再発時間診断指標のプロット

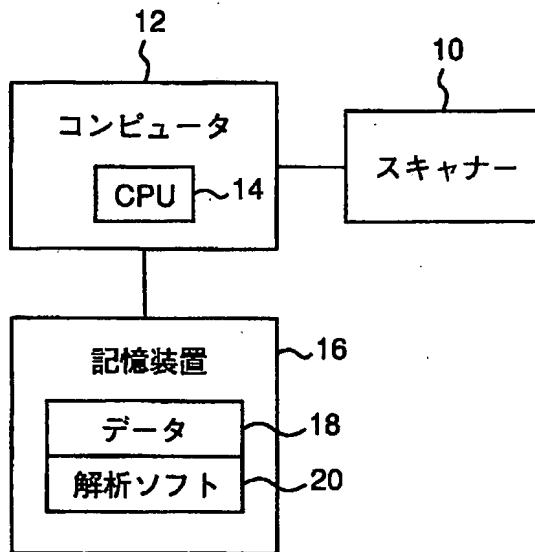
【図 1 5】 実施例 7 の乳癌患者の変数削除基準として $P \geq 0.025$ を採用したときの再発時間診断指標のプロット

【符号の説明】

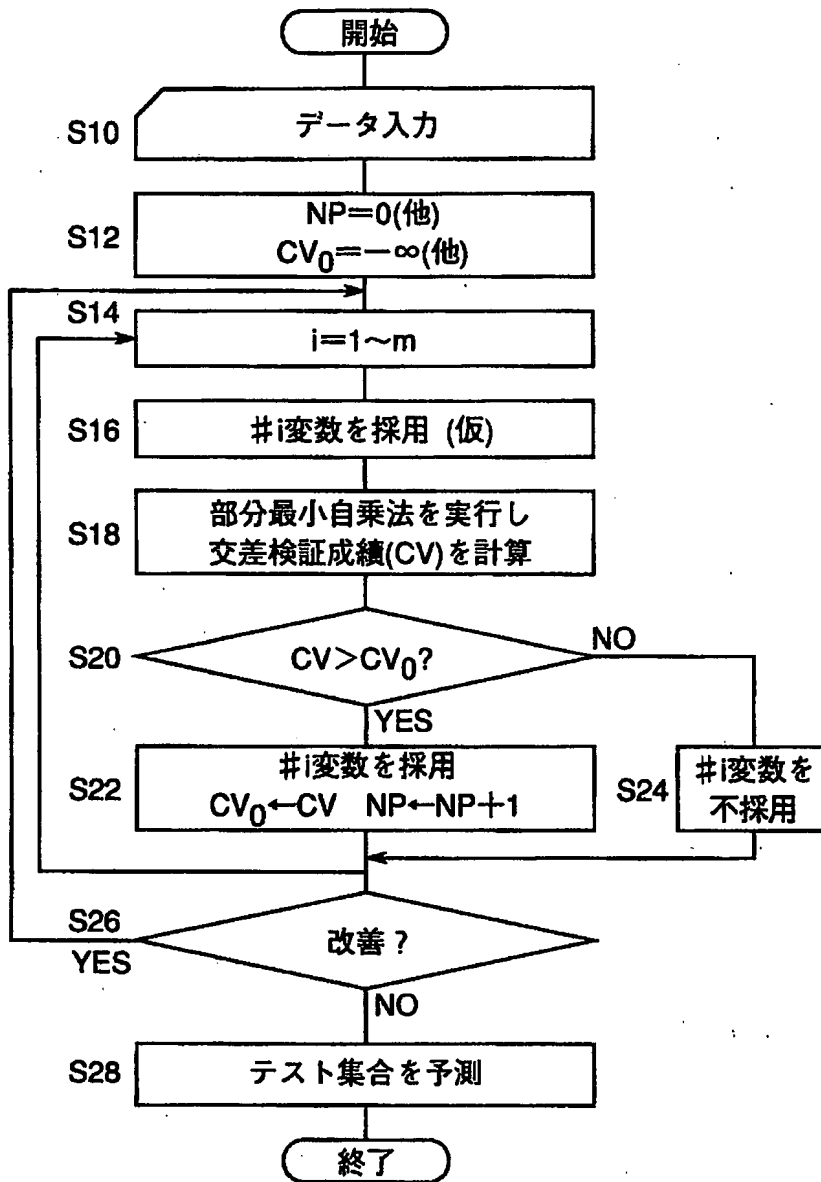
1 0 スキャナ、 1 2 コンピュータ、 1 4 CPU、 1 6 記憶装置、 1 8 記録媒体、 2 0 解析ソフト。

【書類名】 図面

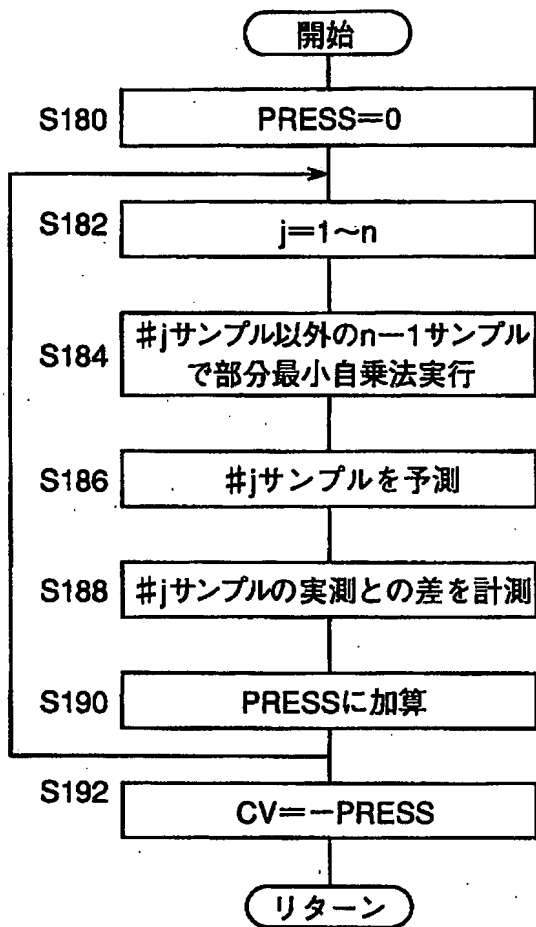
【図1】



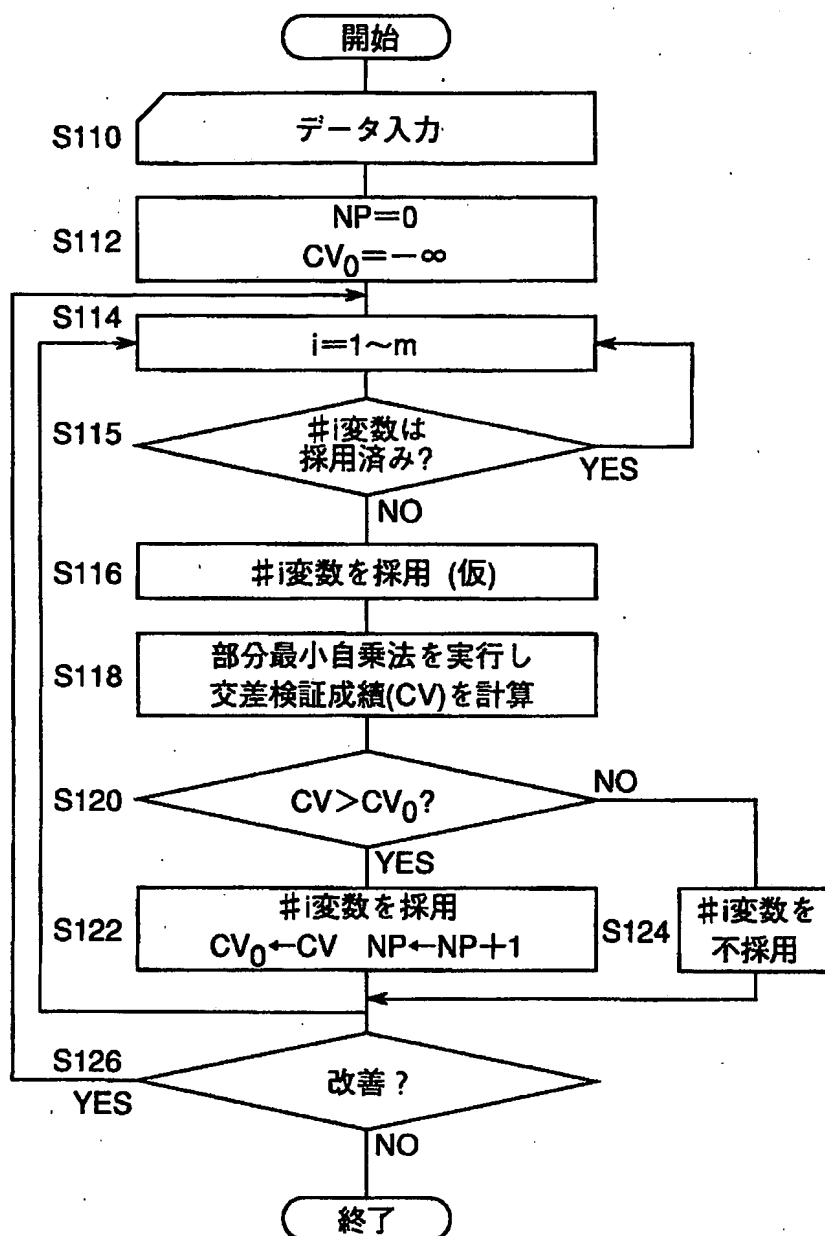
【図 2】



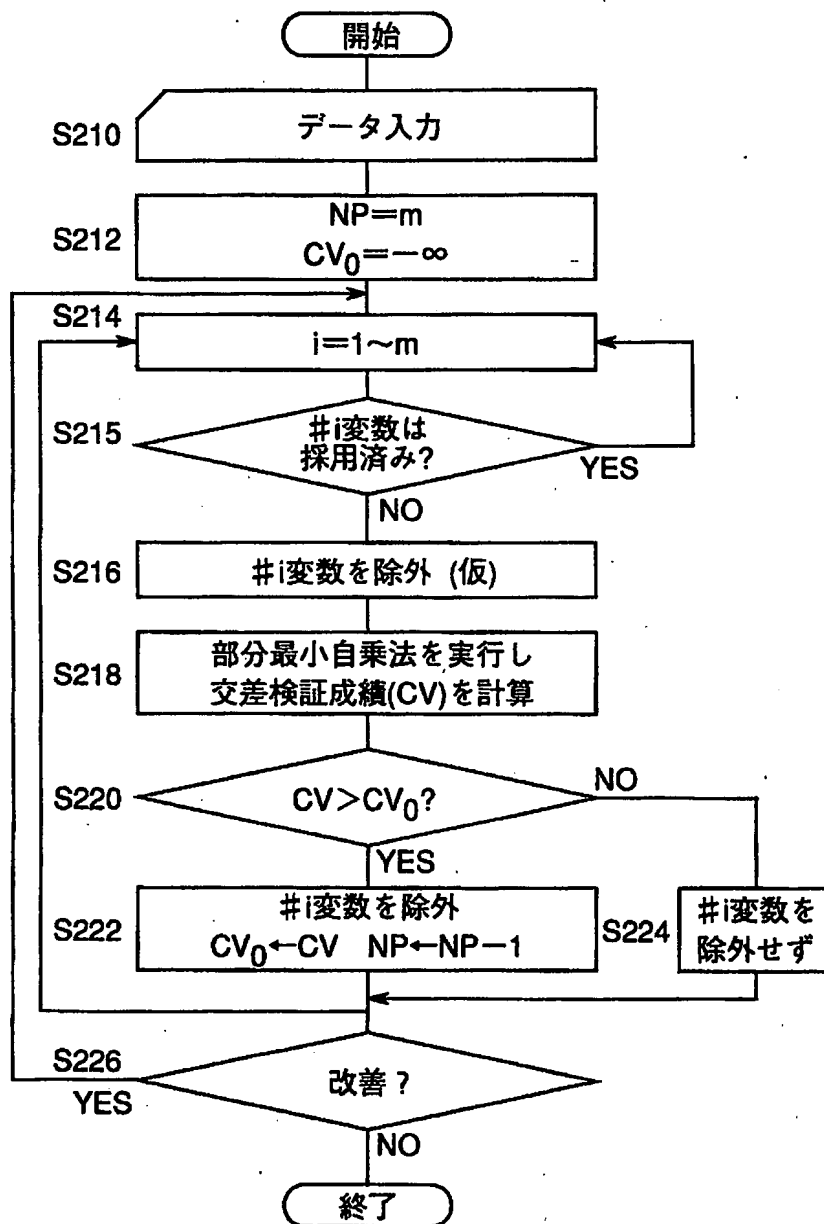
【図 3】



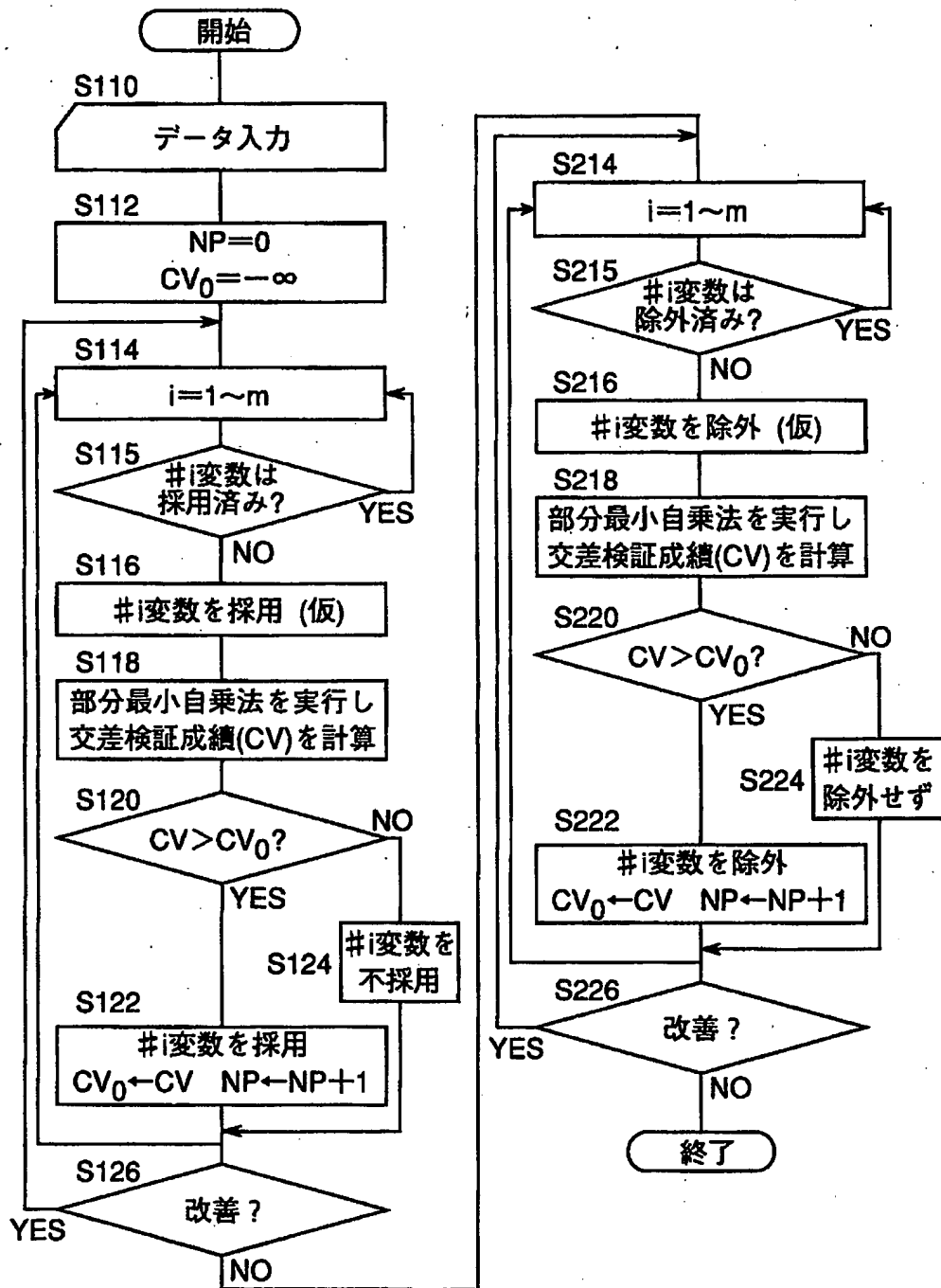
【図4】



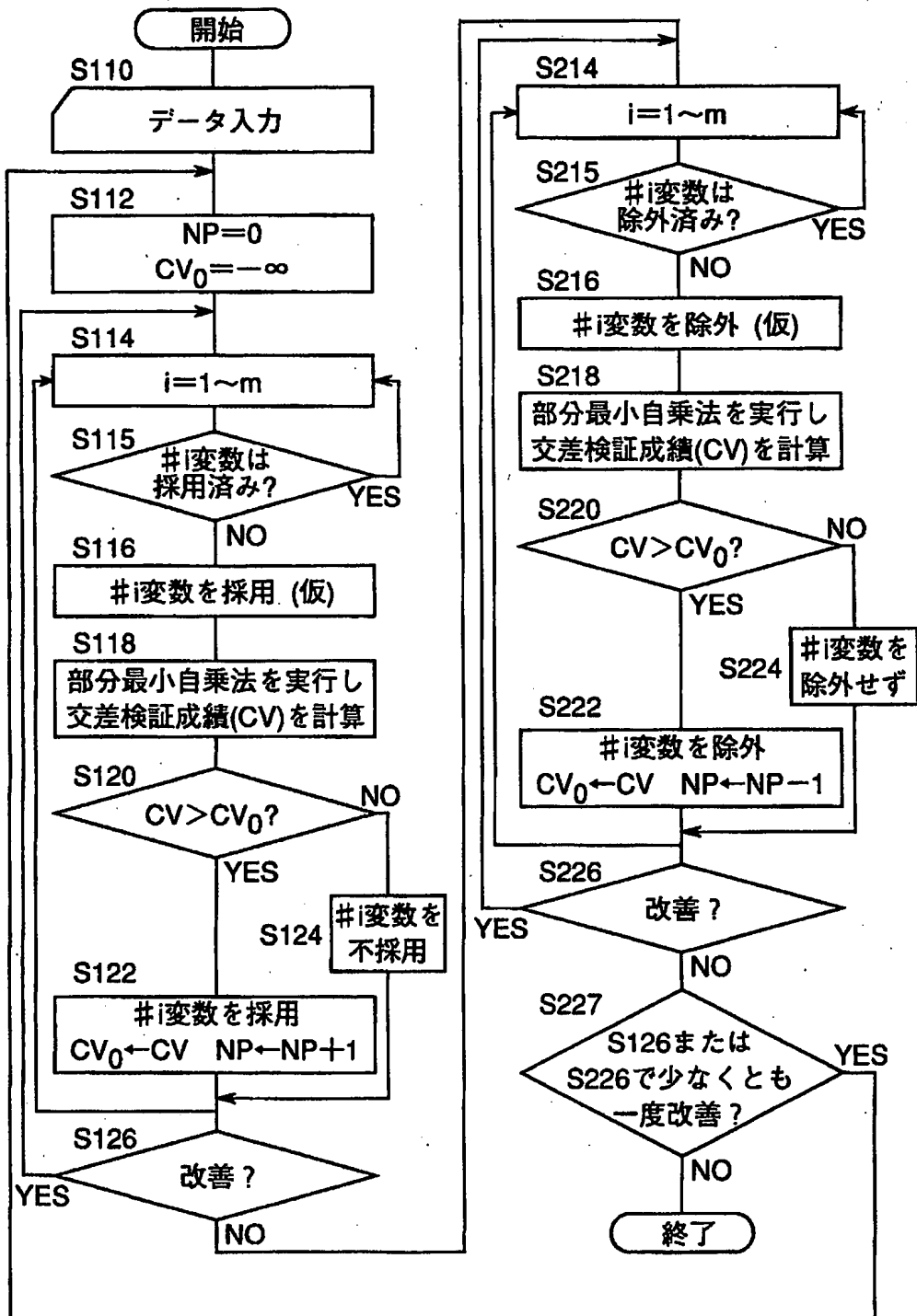
【図 5】



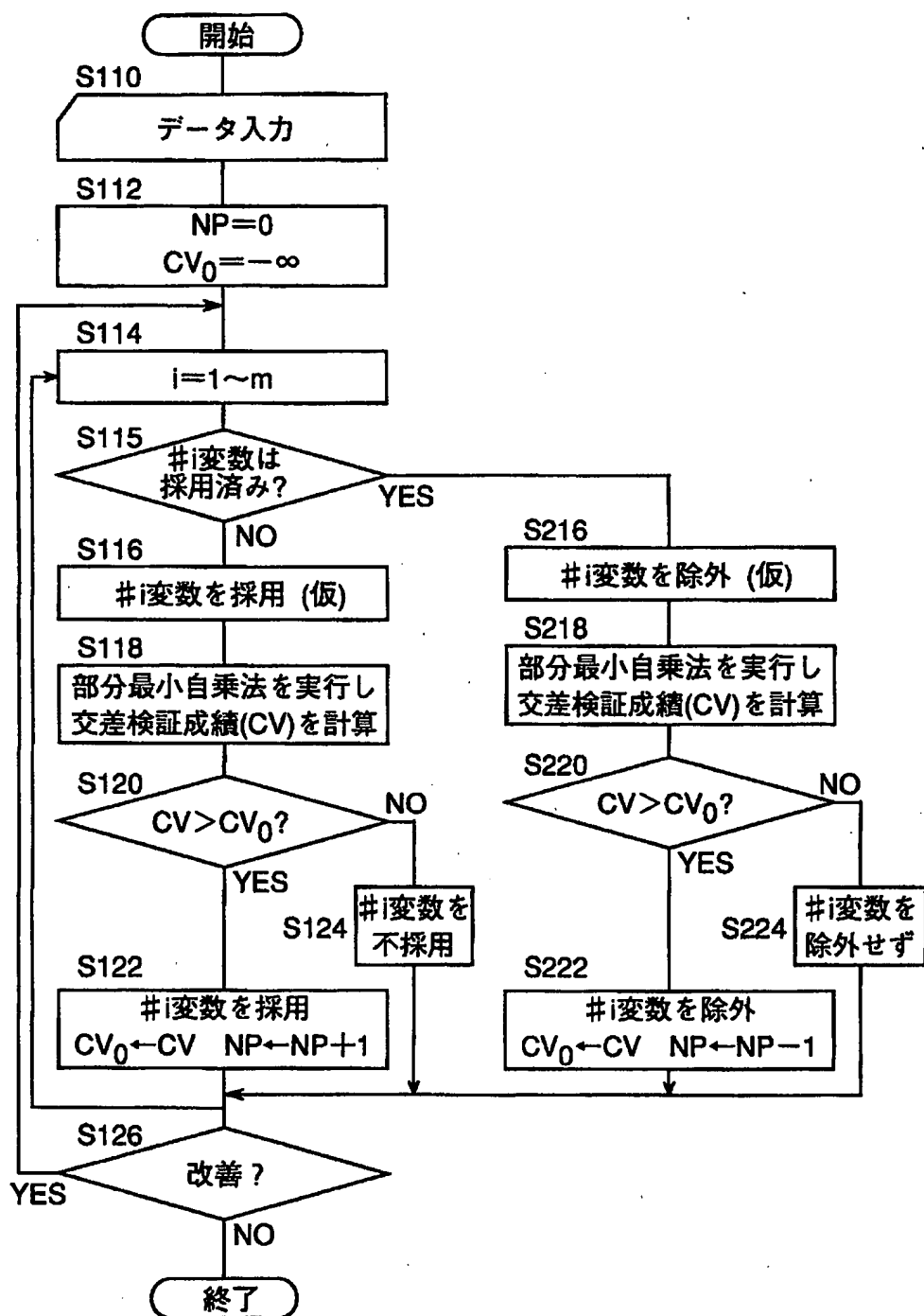
【図 6】



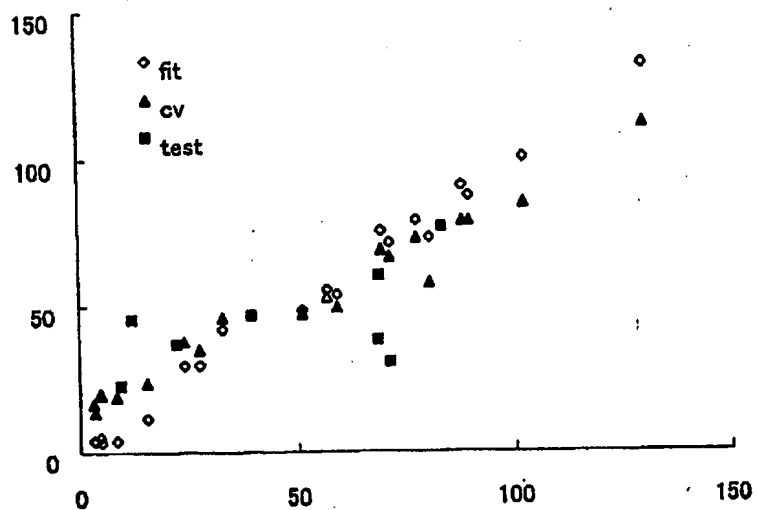
【図7】



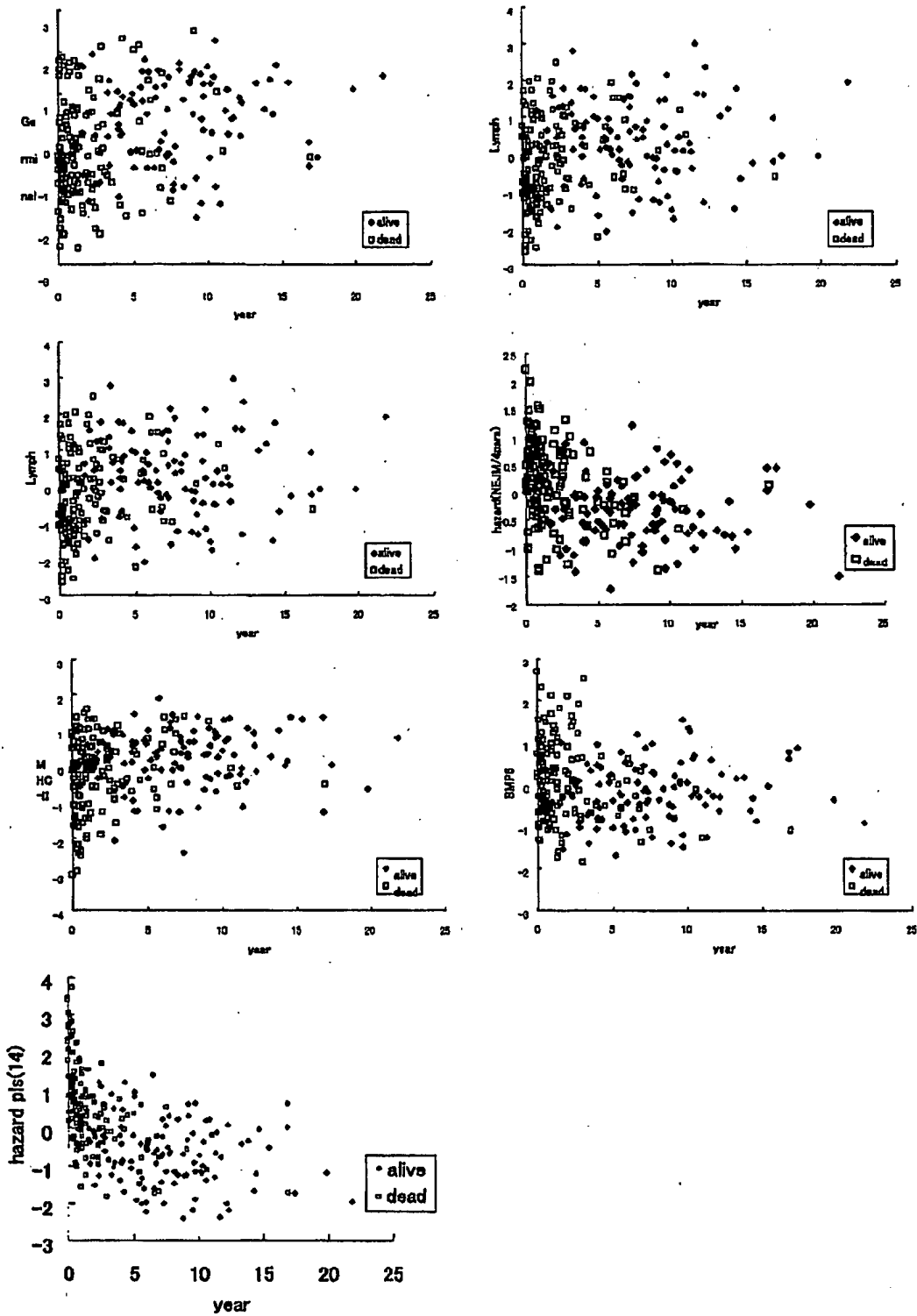
【図 8】



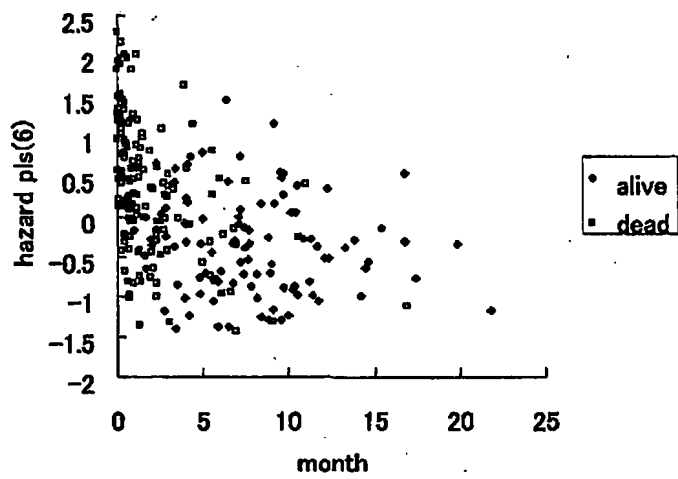
【図9】



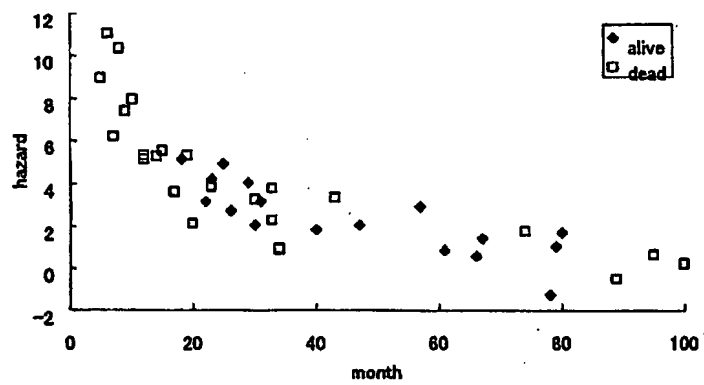
【図10】



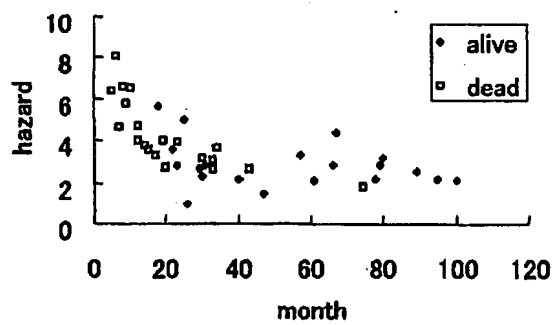
【図11】



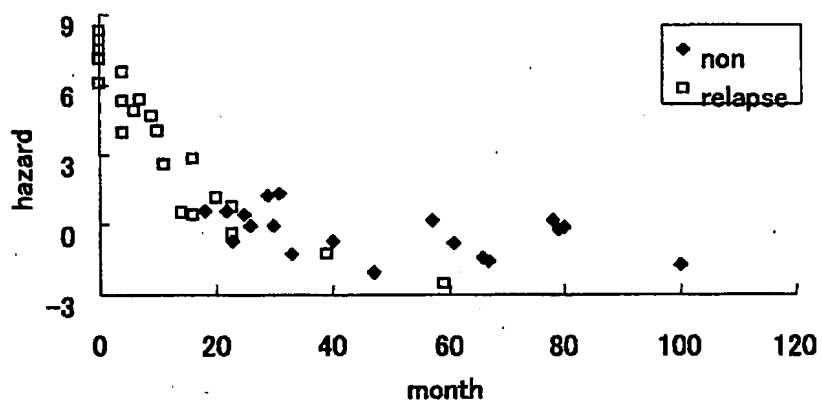
【図12】



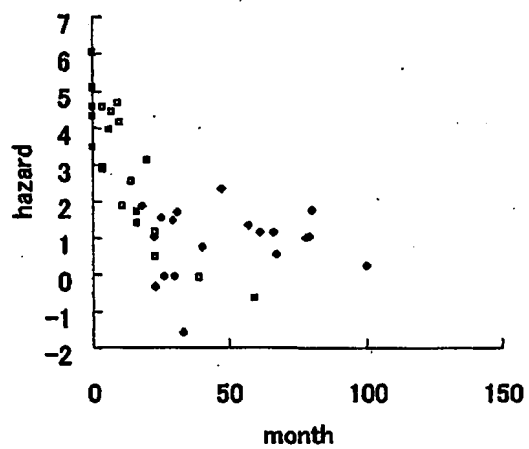
【図13】



【図14】



【図15】



【書類名】 要約書

【要約】

【課題】 多変量の遺伝子発現情報の効果的な情報処理を提供する。

【解決手段】 生体の状態と複数の遺伝子発現の量および／または細胞内物質の量との相関モデルを決定するデータ解析において、

生体の状態を目的変数とし、複数(m 個)の遺伝子発現の量および／または細胞内物質の量を説明変数とするデータの集合において、データに含まれる説明変数を選択し、選択された説明変数と目的変数とを含む相関モデルについて交差検証成績を計算し、その結果を評価判定する。ここで、交差検証成績が改善しなくなるまで、説明変数の選択、交差検証成績の計算、その結果の評価判定を行い、部分最小自乗法モデルを決定する。

【選択図】 図 2

特2002-352645

出願人履歴情報

識別番号 [000000354]

1. 変更年月日 1993年 6月21日

[変更理由] 住所変更

住 所 大阪府大阪市西区江戸堀一丁目3番15号
氏 名 石原産業株式会社